

# BEST AVAILABLE COPY



Europäisches Patentamt  
European Patent Office  
Office européen des brevets



⑪ Publication number: **0 570 729 A2**

⑫

## EUROPEAN PATENT APPLICATION

⑮ Application number: 93106814.2

⑯ Int. Cl.<sup>5</sup> G06F 15/16

⑰ Date of filing: 27.04.93

⑱ Priority: 22.05.92 US 887258

⑲ Date of publication of application:  
24.11.93 Bulletin 93/47

⑳ Designated Contracting States:  
DE FR GB

㉑ Applicant: **INTERNATIONAL BUSINESS  
MACHINES CORPORATION**  
Old Orchard Road  
Armonk, N.Y. 10504(US)

㉒ Inventor: Collins, Olive Allen  
8 Monroe Drive  
Poughkeepsie, New York 12601(US)  
Inventor: Dapp, Michael Charles  
1130 Iyon Avenue  
Endwell, New York 13760(US)  
Inventor: Dieffenderfer, James Warren  
398 Front Street  
Owego, New York 13827(US)  
Inventor: Kuchinski, David Christopher  
19 Gail Drive  
Owego, New York 13827(US)  
Inventor: Knowles, Billy Jack

72 Hurley Avenue  
Kingston, New York 12401(US)  
Inventor: Nier, Richard Edward  
198 Forest Hill Road  
Apalachin, New York 13820(US)  
Inventor: Retter, Eric Eugene  
HCR 34, Box 29B  
Warren Center, Pennsylvania 13827(US)  
Inventor: Richardson, Robert Reist  
RD No.2, Mason Road,  
Box 81  
Vestal, New York 13827(US)  
Inventor: Rolfe, David Bruce  
24 Pine Tree Road  
West Hurley, New York 12491(US)  
Inventor: Smoral, Vincent John  
812 Skyland Terrace  
Endwell, New York 13760(US)

㉓ Representative: Schäfer, Wolfgang, Dipl.-Ing.  
IBM Deutschland Informationssysteme  
GmbH,  
Patentwesen und Urheberrecht  
D-70548 Stuttgart (DE)

㉔ Aapsp I/O programmable router.

EP 0 570 729 A2

㉕ A parallel array processor for massively parallel applications is formed with low power CMOS with DRAM processing while incorporating processing elements on a single chip. Eight processors on a single chip have their own associated processing element, significant memory, and I/O and are interconnected with a hypercube based, but modified, topology. These nodes are then interconnected, either by a hypercube, modified hypercube, or ring, or ring within ring network topology. Conventional microprocessor MIPs consume pins and time going to memory. The new architecture merges processor and memory with multiple PMEs (eight 16 bit processors with 32K and I/O) in DRAM and has no memory access delays and uses all the pins for networking. The chip can be a single node of a fine-grained parallel processor. Each chip will have eight 16 bit processors, each processor providing 5 MIPs performance. I/O has three internal ports and one external port shared by the plural processors on the chip. Significant software flexibility is provided to enable quick implementation of existing programs written in common languages. It is a developable and expandable technology without need to develop new pinouts, new software, or new utilities as chip density increases and new hardware is provided for a chip function. The scalable chip PME has internal and external connections for broadcast and asynchronous SIMD, MIMD and SIMIMD (SIMD/MIMD) with dynamic switching of modes. The chip can be used in systems which employ 32, 64 or 128,000 processors, and can be used for lower, intermediate and higher ranges. Local



FIELD OF THE INVENTIONS

The invention relates to computer and computer systems and particularly to parallel array processors. In accordance with the invention, a parallel array processor (APAP) may be incorporated on a single semiconductor silicon chip. This chip forms a basis for the systems described which are capable of massively parallel processing of complex scientific and business applications.

REFERENCES USED IN THE DISCUSSION OF THE INVENTIONS

In the detailed discussion of the invention, other works will be referenced, including references to our own unpublished works which are not Prior Art, which will aid the reader in following the discussion.

GLOSSARY OF TERMS

- o ALU  
ALU is the arithmetic logic unit portion of a processor.
- o Array  
Array refers to an arrangement of elements in one or more dimensions. An array can include an ordered set of data items (array element) which in languages like Fortran are identified by a single name. In other languages such a name of an ordered set of data items refers to an ordered collection or set of data elements, all of which have identical attributes. A program array has dimensions specified, generally by a number or dimension attribute. The declarator of the array may also specify the size of each dimension of the array in some languages. In some languages, an array is an arrangement of elements in a table. In a hardware sense, an array is a collection of structures (functional elements) which are generally identical in a massively parallel architecture. Array elements in data parallel computing are elements which can be assigned operations and when parallel can each independently and in parallel execute the operations required. Generally, arrays may be thought of as grids of processing elements. Sections of the array may be assigned sectional data, so that sectional data can be moved around in a regular grid pattern. However, data can be indexed or assigned to an arbitrary location in an array.
- o Array Director  
An Array Director is a unit programmed as a controller for an array. It performs the function of a master controller for a grouping of functional elements arranged in an array.
- o Array Processor  
There two principal types of array processors - multiple instruction multiple data (MIMD) and single instruction multiple data (SIMD). In a MIMD array processor, each processing element in the array executes its own unique instruction stream with its own data. In a SIMD array processor, each processing element in the array is restricted to the same instruction via a common instruction stream; however, the data associated with each processing element is unique. Our preferred array processor has other characteristics. We call it Advanced Parallel Array Processor, and use the acronym APAP.
- o Asynchronous  
Asynchronous is without a regular time relationship; the execution of a function is unpredictable with respect to the execution of other functions which occur without a regular or predictable time relationship to other function executions. In control situations, a controller will address a location to which control is passed when data is waiting for an idle element being addressed. This permits operations to remain in a sequence while they are out of time coincidence with any event.
- o BOPS/GOPS  
BOPS or GOPS are acronyms having the same meaning - billions of operations per second. See GOPS.
- o Circuit Switched/Store Forward  
These terms refer to two mechanisms for moving data packets through a network of nodes. Store Forward is a mechanism whereby a data packet is received by each intermediate node, stored into its memory, and then forwarded on towards its destination. Circuit Switch is a mechanism whereby an intermediate node is commanded to logically connect its input port to an output port such that data packets can pass directly through the node towards their destination, without entering the intermediate node's memory.
- o Cluster  
A cluster is a station (or functional unit) which consists of a control unit (cluster controller) and the

hardware (which may be terminals, functional units, or virtual components) attached to it. Our Cluster includes an array of PMEs sometimes called a Node array. Usually a cluster has 512 PMEs.

Our Entire PME node array consists of a set of clusters, each cluster supported by a cluster controller (CC).

5   o Cluster controller

A cluster controller is a device that controls input/output (I/O) operations for more than one device or functional unit connected to it. A cluster controller is usually controlled by a program stored and executed in the unit as it was in the IBM 3601 Finance Communication Controller, but it can be entirely controlled by hardware as it was in the IBM 3272 Control Unit.

10   o Cluster synchronizer

A cluster synchronizer is a functional unit which manages the operations of all or part of a cluster to maintain synchronous operation of the elements so that the functional units maintain a particular time relationship with the execution of a program.

15   o Controller

A controller is a device that directs the transmission of data and instructions over the links of an interconnection network; its operation is controlled by a program executed by a processor to which the controller is connected or by a program executed within the device.

20   o CMOS

CMOS is an acronym for Complementary Metal-Oxide Semiconductor technology. It is commonly used to manufacture dynamic random access memories (DRAMs). NMOS is another technology used to manufacture DRAMS. We prefer CMOS but the technology used to manufacture the APAP is not intended to limit the scope of the semiconductor technology which is employed.

25   o Dotting

Dotting refers to the joining of three or more leads by physically connecting them together. Most backpanel busses share this connection approach. The term relates to OR DOTS of times past but is used here to identify multiple data sources that can be combined onto a bus by a very simple protocol.

Our I/O zipper concept can be used to implement the concept that the port into a node could be driven by the port out of a node or by data coming from the system bus. Conversely, data being put out of a node would be available to both the input to another node and to the system bus. Note that outputting data to both the system bus and another node is not done simultaneously but in different cycles.

Dotting is used in the H-DOT discussions where Two-ported PEs or PMEs or Pickels can be used in arrays of various organizations by taking advantage of dotting. Several topologies are discussed including 2D and 3D Meshes, Base 2 N-cube, Sparse Base 4 N-cube, and Sparse Base 8 N-cube.

30   o DRAM

DRAM is an acronym for dynamic random access memory, the common storage used by computers for main memory. However, the term DRAM can be applied to use as a cache or as a memory which is not the main memory.

35   o FLOATING-POINT

A floating-point number is expressed in two parts. There is a fixed point or fraction part, and an exponent part to some assumed radix or Base. The exponent indicates the actual placement of the decimal point. In the typical floating-point representation a real number 0.0001234 is represented as 0.1234-3, where 0.1234 is the fixed-point part and -3 is the exponent. In this example, the floating-point radix or base is 10, where 10 represents the implicit fixed positive integer base, greater than unity, that is raised to the power explicitly denoted by the exponent in the floating-point representation or represented by the characteristic in the floating-point representation and then multiplied by the fixed-point part to determine the real number represented. Numeric literals can be expressed in floating-point notation as well as real numbers.

40   o FLOPS

This terms refers to floating-point instructions per second. Floating-point operations include ADD, SUB, MPY, DIV and often many others. Floating-point instructions per second parameter is often calculated using the add or multiply instructions and, in general, may be considered to have a 50/50 mix. An operation includes the generation of exponent, fraction and any required fraction normalization. We could address 32 or 48-bit floating-point formats (or longer but we have not counted them in the mix.) A floating-point operation when implemented with fixed point instructions (normal or RISC) requires multiple instructions. Some use a 10 to 1 ratio in figuring performance while some specific studies have shown a ratio of 8.25 more appropriate to use. Various architectures will have different

EP 0 570 729 A2

ratios.

o Functional unit

A functional unit is an entity of hardware, software, or both, capable of accomplishing a purpose.

o Gbytes

5 Gbytes refers to a billion bytes. Gbytes/s would be a billion bytes per second.

o GIGAFLOPS

(10)<sup>9</sup> floating-point instructions per second.

o GOPS and PETAOPS

70 GOPS or BOPS, have the same meaning - billions of operations per second. PETAOPS means trillions of operations per second, a potential of the current machine. For our APAP machine they are just about the same as BIPs/GIPs meaning billions of instructions per second. In some machines an instruction can cause two or more operations (ie. both an add and multiply) but we don't do that. Alternatively it could take many instructions to do an op. For example we use multiple instructions to perform 64 bit arithmetic. In counting ops however, we did not elect to count log ops. GOPS may be the preferred use to describe performance, but there is no consistency in usage that has been noted. One sees MIPs/MOPs then BIPs/BOPs and MegaFLOPS/GigaFLOPS/TeraFLOPS/PetaFlops.

o ISA

ISA means the Instruction Set Architecture.

o Link

20 A link is an element which may be physical or logical. A physical link is the physical connection for joining elements or units, while in computer programming a link is an instruction or address that passes control and parameters between separate portions of the program. In multisystems a link is the connection between two systems which may be specified by program code identifying the link which may be identified by a real or virtual address. Thus generally a link includes the physical medium, any protocol, and associated devices and programming; it is both logical and physical.

o MFLOPS

MFLOPS means (10)<sup>6</sup> floating-point instructions per second.

o MIMD

30 MIMD is used to refer to a processor array architecture wherein each processor in the array has its own instruction stream, thus Multiple Instruction stream, to execute Multiple Data streams located one per processing element.

o Module

35 A module is a program unit that is discrete and identifiable or a functional unit of hardware designed for use with other components. Also, a collection of PEs contained in a single electronic chip is called a module.

o Node

Generally, a node is the junction of links. In a generic array of PEs, one PE can be a node. A node can also contain a collection of PEs called a module. In accordance with our invention a node is formed of an array of PMEs, and we refer to the set of PMEs as a node. Preferably a node is 8 PMEs.

40 o Node array

A collection of modules made up of PMEs is sometimes referred to as a node array, is an array of nodes made up of modules. A node array is usually more than a few PMEs, but the term encompasses a plurality.

o PDE

45 A PDE is a partial differential equation.

o PDE relaxation solution process

PDE relaxation solution process is a way to solve a PDE (partial differential equation). Solving PDEs uses most of the super computing compute power in the known universe and can therefore be a good example of the relaxation process. There are many ways to solve the PDE equation and more than one of the numerical methods includes the relaxation process. For example, if a PDE is solved by finite element methods relaxation consumes the bulk of the computing time. Consider an example from the world of heat transfer. Given hot gas inside a chimney and a cold wind outside, how will the temperature gradient within the chimney bricks develop? By considering the bricks as tiny segments and writing an equation that says how heat flows between segments as a function of temperature differences then the heat transfer PDE has been converted into a finite element problem. If we then say all elements except those on the inside and outside are at room temperature while the boundary segments are at the hot gas and cold wind temperature, we have set up the problem to begin relaxation. The computer program then models time by updating the temperature variable in each

segment based upon the amount of heat that flows into or out of the segment. It takes many cycles of processing all the segments in the model before the set of temperature variables across the chimney relaxes to represent actual temperature distribution that would occur in the physical chimney. If the objective was to model gas cooling in the chimney then the elements would have to extend to gas equations, and the boundary conditions on the inside would be linked to another finite element model, and the process continues. Note that the heat flow is dependent upon the temperature difference between the segment and its neighbors. It thus uses the inter-PE communication paths to distribute the temperature variables. It is this near neighbor communication pattern or characteristic that makes PDE relation very applicable to parallel computing.

#### 10 ○ PICKET

This is the element in an array of elements making up an array processor. It consists of: data flow (ALU REGS), memory, control, and the portion of the communication matrix associated with the element. The unit refers to a 1/nth of an array processor made up of parallel processor and memory elements with their control and portion of the array intercommunication mechanism. A picket is a form of processor memory element or PME. Our PME chip design processor logic can implement the picket logic described in related applications or have the logic for the array of processors formed as a node. The term PICKET is similar to the commonly used array term PE for processing element, and is an element of the processing array preferably comprised of a combined processing element and local memory for processing bit parallel bytes of information in a clock cycle. The preferred embodiment consisting of a byte wide data flow processor, 32k bytes or more of memory, primitive controls and ties to communications with other pickets.

The term "picket" comes from Tom Sawyer and his white fence, although it will also be understood functionally that a military picket line analogy fits quite well.

#### 20 ○ Picket Chip

A picket chip contains a plurality of pickets on a single silicon chip.

#### 25 ○ Picket Processor system (or Subsystem)

A picket processor is a total system consisting of an array of pickets, a communication network, an I/O system, and a SIMD controller consisting of a microprocessor, a canned routine processor, and a micro-controller that runs the array.

#### 30 ○ Picket Architecture

The Picket Architecture is the preferred embodiment for the SIMD architecture with features that accommodate several diverse kinds of problems including:

- set associative processing
- parallel numerically intensive processing
- 35 - physical array processing similar to images

#### 40 ○ Picket Array

A picket array is a collection of pickets arranged in a geometric order, a regular array.

#### 45 ○ PME or processor memory element

PME is used for a processor memory element. We use the term PME to refer to a single processor, memory and I/O capable system element or unit that forms one of our parallel array processors. A processor memory element is a term which encompasses a picket. A processor memory element is 1/nth of a processor array which comprises a processor, its associated memory, control interface, and a portion of an array communication network mechanism. This element can have a processor memory element with a connectivity of a regular array, as in a picket processor, or as part of a subarray, as in the multi-processor memory element node we have described.

#### 50 ○ Routing

Routing is the assignment of a physical path by which a message will reach its destination. Routing assignments have a source or origin and a destination. These elements or addresses have a temporary relationship or affinity. Often, message routing is based upon a key which is obtained by reference to a table of assignments. In a network, a destination is any station or network addressable unit addressed as the destination of information transmitted by a path control address that identifies the link. The destination field identifies the destination with a message header destination code.

#### 55 ○ SIMD

A processor array architecture wherein all processors in the array are commanded from a Single Instruction stream to execute Multiple Data streams located one per processing element.

#### 60 ○ SIMD/MIMD or SIMD/MIMD

SIMD/MIMD or SIMD/MIMD is a term referring to a machine that has a dual function that can switch from MIMD to SIMD for a period of time to handle some complex instruction, and thus has two

EP 0 570 729 A2

modes. The Thinking Machines, Inc. Connection Machine model CM-2 when placed as a front end or back end of a MIMD machine permitted programmers to operate different modes for execution of different parts of a problem, referred to sometimes as dual modes. These machines have existed since Illiac and have employed a bus that interconnects the master CPU with other processors. The master control processor would have the capability of interrupting the processing of other CPUs. The other CPUs could run independent program code. During an interruption, some provision must be made for checkpointing (closing and saving current status of the controlled processors).

#### o SIMIMD

SIMIMD is a processor array architecture wherein all processors in the array are commanded from a Single Instruction stream, to execute Multiple Data streams located one per processing element. Within this construct, data dependent operations within each picket that mimic instruction execution are controlled by the SIMD instruction stream.

This is a Single Instruction Stream machine with the ability to sequence Multiple Instruction streams (one per Picket) using the SIMD instruction stream and operate on Multiple Data Streams (one per Picket). SIMIMD can be executed by a processor memory element system.

#### SISD

SISD is an acronym for Single Instruction Single Data.

#### o Swapping

Swapping interchanges the data content of a storage area with that of another area of storage.

#### o Synchronous Operation

Synchronous operation in a MIMD machine is a mode of operation in which each action is related to an event (usually a clock); it can be a specified event that occurs regularly in a program sequence. An operation is dispatched to a number of PEs who then go off to independently perform the function. Control is not returned to the controller until the operation is completed. If the request is to an array of functional units, the request is generated by a controller to elements in the array which must complete their operation before control is returned to the controller.

#### o TERAFLIPS

TERAFLIPS means  $(10)^{12}$  floating-point instructions per second.

#### o VLSI

VLSI is an acronym for very large scale integration (as applied to integrated circuits).

#### o Zipper

A zipper is a new function provided. It allows for links to be made from devices which are external to the normal interconnection of an array configuration.

#### BACKGROUND OF THE INVENTION

In the never ending quest for faster computers, engineers are linking hundreds, and even thousands of low cost microprocessors together in parallel to create super supercomputers that divide in order to conquer complex problems that stump today's machines. Such machines are called massively parallel. We have created a new way to create massively parallel systems. The many improvements which we have made should be considered against the background of many works of others.

Multiple computers operating in parallel have existed for decades. Early parallel machines included the ILLIAC which was started in the 1960s. ILLIAC IV was built in the 1970s. Other multiple processors include (see a partial summary in U.S. Patent 4,975,834 Issued December 4, 1990 to Xu et al) the Cedar, Sigma-1, the Butterfly and the Monarch, the Intel ipso, The Connection Machines, the Caltech COSMIC, the N Cube, IBM's RP3, IBM's GF11, the NYU Ultra Computer, the Intel Delta and Touchstone.

Large multiple processors beginning with ILLIAC have been considered supercomputers. Supercomputers with greatest commercial success have been based upon multiple vector processors, represented by the Cray Research Y-MP systems, the IBM 3090, and other manufacturer's machines including those of Amdahl, Hitachi, Fujitsu, and NEC.

Massively Parallel Processors (MPPs) are now thought of as capable of becoming supercomputers. These computer systems aggregate a large number of microprocessors with an interconnection network and program them to operate in parallel. There have been two modes of operation of these computers. Some of these machines have been MIMD mode machines.

Some of these machines have been SIMD mode machines. Perhaps the most commercially acclaimed of these machines has been the Connection Machines series 1 and 2 of Thinking Machines, Inc. These

have been essentially SIMD machines. Many of the massively parallel machines have used microprocessors interconnected in parallel to obtain their concurrency or parallel operations capability. Intel microprocessors like i860 have been used by Intel and others. nCube has made such machines with Intel '386 microprocessors. Other machines have been built with what is called the "transputer" chip. Immos Transputer IMS T800 is an example. The Immos Transputer T800 is a 32 bit device with an integral high speed floating point processor.

As an example of the kind of systems that are built, several Immos Transputer T800 chips each would have 32 communication link inputs and 32 link outputs. Each chip would have a single processor, a small amount of memory, and communication links to the local memory and to an external interface. In addition, in order to build up the system communication link adaptors like IMS C011 and C012 would be connected. In addition switches, like a IMS C004 would provide, say, a crossbar switch between the 32 link inputs and 32 link outputs to provide point-to-point connection between additional transputer chips. In addition, there will be special circuitry and interface chips for transputers adapting them to be used for a special purpose tailored to the requirements of a specific device, a graphics or disk controller. The Immos IMS M212 is a 18 bit processor, with on chip memory and communication links. It contains hardware and logic to control disk drives and can be used as a programmable disk controller or as a general purpose interface. In order to use the concurrency (parallel operations) Immos developed a special language, Occam, for the transputer. Programmers have to describe the network of transputers directly in an Occam program.

Some of these massively parallel machines use parallel processor arrays of processor chips which are interconnected with different topologies. The transputer provides a crossbar network with the addition of IMS C004 chips. Some other systems use a hypercube connection. Others use a bus or mesh to connect the microprocessors and their associated circuitry. Some have been interconnected by circuit switch processors that use switches as processor addressable networks. Generally, as with the 14 RISC/6000s which were interconnected last fall at Lawrence Livermore by wiring the machines together, the processor addressable networks have been considered as coarse-grained multiprocessors.

Some very large machines are being built by Intel and nCube and others to attack what are called "grand challenges" in data processing. However, these computers are very expensive. Recent projected costs are in the order of \$30,000,000.00 to \$75,000,000.00 (Tera Computer) for computers whose development has been funded by the U.S. Government to attack the "grand challenges". These "grand challenges" would include such problems as climate modeling, fluid turbulence, pollution dispersion, mapping of the human genome and ocean circulation, quantum chromodynamics, semiconductor and supercomputer modeling, combustion systems, vision and cognition.

As a footnote to our background, we should recognize one of the early massively parallel machines developed by IBM. In our description we have chosen to use the term processor memory element rather than "transputer" to describe one of the eight or more memory units with processor and I/O capabilities which make up the array of PME's in a chip, or node. The referenced prior art "transputer" has on a chip one processor, a Fortran coprocessor, and a small memory, with an I/O interface. Our processor memory element could apply to a transputer and to the PME of the RP3 generally. However, as will be recognized, our little chip is significantly different in many respects. Our little chip has many features described later. However, we do recognize that the term PME was first coined for another, now more typical, PME which formed the basis for the massively parallel machine known as the RP3. The IBM Research Parallel Processing Prototype (RP3) was an experimental parallel processor based on a Multiple Instruction Multiple Data (MIMD) architecture. RP3 was designed and built at IBM T.J. Watson Research Center in cooperation with the New York University Ultracomputer project. This work was sponsored in part by Defense Advanced Research Project Agency. RP3 was comprised of 64 Processor-Memory Elements (PMEs) interconnected by a high speed omega network. Each PME contained a 32-bit IBM "PC scientific" microprocessor, 32-kB cache, a 4-MB segment of the system memory, and an I/O port. The PME I/O port hardware and software supported initialization, status acquisition, as well as memory and processor communication through shared I/O support Processors (ISPs). Each ISP supports eight processor-memory elements through the Extended I/O adapters (ETIOs), independent of the system network. Each ISP interfaced to the IBM S/370 channel and the IBM Token-Ring network as well as providing operator monitor service. Each extended I/O adapter attached as a device to a PME ROMP Storage Channel (RSC) and provided programmable PME control/status signal I/O via the ETIO channel. The ETIO channel is the 32-bit bus which interconnected the ISP to the eight adapters. The ETIO channel relied on a custom interface protocol which was supported by hardware on the ETIO adapter and software on the ISP.



EP 0 570 729 A2

Problems addressed by our APAP machine

The machine which we have called the Advanced Parallel Array Processor (APAP) is a fine-grained parallel processor which we believe is needed to address issues of prior designs. As illustrated above, there have been many fine-grained (and also coarse-grained) processors constructed from both point design and off-the-shelf processors using dedicated and shared memory and any one of the many possible interconnection schemes. To date these approaches have all encountered one or more design and performance limitations. Each "solution" leads in a different direction. Each has its problems. Existing parallel machines are difficult to program. Each is not generally adaptable to various sizes of machines compatible across a range of applications. Each has its design limitations caused by physical design, interconnection and architectural issues.

Physical Issues

Some approaches utilize a separate chip design for each of the various functions required in a horizontal structure. These approaches suffer performance limitations due to chip crossing delays.

Other approaches integrate various functions together vertically into a single chip. These approaches suffer performance limitations due to the physical limit on the number of logic gates which can be integrated onto a producible chip.

Interconnection Issues

Networks which interconnect the various processing functions are important to fine-grained parallel processors. Processor designs with buses, meshes, and hypercubes have all been developed. Each of these networks has inherent limitations as to processing capability. Buses limit both the number of processors which can be physically interconnected and the network performance. Meshes lead to large network diameters which limit network performance. Hypercubes require each node to have a large number of interconnection ports; the number of processors which can be interconnected is limited by the physical input/output pins at the node. Hypercubes are recognized as having some significant performance gains over the prior bus and mesh structures.

Architectural Issues:

Processes which are suitable for fine-grained parallel processors fall into two distinct types. Processes which are functionally partitionable tend to perform better on multiple instruction, multiple data (MIMD) architectures. Processes which are not functionally partitionable but have multiple data streams tend to perform better on single instruction, multiple data (SIMD) architectures. For any given application, there is likely to be some number of both types of processes. System trade-offs are required to pick the architecture which best suits a particular application but no single solution has been satisfactory.

SUMMARY OF THE INVENTION

We have created a new way to make massively parallel processors and other computer systems by creating a new "chip" and systems designed with our new concepts. This application is directed to such systems. Components described in our applications can be combined in our systems to make new systems. They also can be combined with existing technology.

Think, our little CMOS DRAM chip of approximately 14 x 14 mm can be put together much like bricks are walled in a building or paved to form a brick road. Our chip provides the structure necessary to build a "house", a complex computer system, by connected replication.

Placing our development in perspective, four little chips, each one alike, each one with eight or more processors embedded in memory with an internal array capability and external I/O broadcast and control interface, would provide the memory and processing power of thirty-six or more complex computers, and they could all be placed with compact hybrid packaging into something the size of a watch, and operated with very low power, as each chip only dissipates about 2 watts. With this chip, we have created many new concepts, and those that we consider our invention are described in detail in the description and claims. The systems that can be created with our computer system can range from small devices to massive machines with PETAOP potential.

Our little memory chip array processor we call our Advanced Parallel Array Processor. Though small, it is complex and powerful. A typical cluster will have many chips.

Many aspects and features of invention have been described in this and related applications. These concepts and features of invention improve and are applicable to computer systems which may not employ each invention. We believe our concepts and features will be adopted and used in the next century.

This technical description provides an overview of our Advanced Parallel Array Processor (APAP) representing our new memory concepts and our effort in developing a scalable massively parallel processor (MPP) that is simple (very small number of unique part numbers) and has very high performance. Our processor utilizes in its preferred embodiment a VLSI chip. The chip comprises  $2^n$  PME microcomputers.  $2^n$  represents the maximum number of array dimensionality. The chip further comprises a broadcast and control interface (BCI) and internal and external communication paths between PMEs on the chip among themselves and to the off chip system environment. The preferred chip has 8 PMEs (but we also can provide more) and one BCI. The  $2^n$  PMEs and BCI are considered a node. This node can function in either SIMD or MIMD mode, in dual SIMD/MODE, with asynchronous processing, and with SIMIMD functionality. Since it is scalable, this approach provides a node which can be the main building block for scalable parallel processors of varying size. The microcomputer architecture of the PME provides fully distributed message passing interconnection and control features within each node, or chip. Each node provides multiple parallel microcomputer capability at the chip level, the microprocessor or personal computer level, at a workstation level, at special application levels which may be represented by a vision and/or avionics level, and, when fully extended, to capability at greater levels with powerful GigaFlop performance into the supercomputer range. The simplicity is achieved by the use of a single highly extended DRAM Chip that is replicated into parallel clusters. This keeps the part number count down and allows scaling capability to the cost or performance need, by varying the chip count, then the number of modules, etc.

Our approach enables us to provide a machine with attributes meeting the requirements that drive to a parallel solution in a series of applications. Our methods of parallelization at the sub-chip level serve to keep weight, volume, and recurring and logistic costs down.

Because our different size systems are all based upon a single chip, software tools are common for all size systems. This offers the potential of development software (running on smaller workstation machines) that is interchangeable among all levels (workstation, aerospace, and supercomputer). That advantage means programmers can develop programs on workstations while a production program runs on a much larger machine.

As a result of our well balanced design implementation we meet today's requirements imposed by technology, performance, cost, and perception, and enable growth of the system into the future. Since our MPP approach starts at the chip level, our discussion starts at the chip technology description and concludes with the supercomputer application descriptions.

Physical, interconnection, and architectural issues will all be addressed in the machine directly. Functions will not only be integrated into a single chip design, but the chip design will provide functions sufficiently powerful and flexible that the chip will be effective at processing, routing, storage and three classes of I/O. The interconnection network will be a new version of the hypercube which provides minimum network diameters without the input/output pin and wireability limitations normally associated with hypercubes. The trade-off between SIMD and MIMD are eliminated because the design allows processors to dynamically switch between MIMD and SIMD mode. This eliminates many problems which will be encountered by application programmers of "hybrid" machines. In addition, the design will allow a subset of the processors to be in SIMD or MIMD mode.

The Advanced Parallel Array Processor (APAP) is a fine-grained parallel processor. It consists of control and processing sections which are partitionable such that configurations suitable for supercomputing through personal computing applications can be satisfied. In most configurations it would attach to a host processor and support the off loading of segments of the host's workload. Because the APAP array processing elements are general purpose computers, the particular type of workload off-loaded will vary depending upon the capabilities of the host. For example, our APAP can be a module for an IBM 3090 vector processor mainframe. When attached to a mainframe with high performance vector floating point capability the task off-loaded might be sparse to dense matrix transformations. Alternatively, when attached to a PC personal computer the off-loaded task might be numerically intensive 3 dimensional graphics processing.

The above referenced parent USN 07/611,584, filed November 13, 1980 of Dieffenenderfer et al., titled "Parallel Associative Processor System" describes the idea of integrating computer memory and control logic within a single chip and replicating the combination within the chip and building a processor system out of replications of the single chip. This approach which is continued and expanded here leads to a

EP 0 570 729 A2

system which provides massively parallel processing capability at the cost of developing and manufacturing only a single chip type while enhancing performance capability by reducing the chip boundary crossings and line length.

The above referenced parent US8N 07/911,594, filed November 13, 1990 illustrated utilization of 1-dimensional I/O structures (essentially a linear I/O) with multiple SIMD PMEs attached to that structure within a chip. This embodiment elaborates these concepts to dimensions greater than 1. The description which follows will be in terms of 4-dimensional I/O structures with 8 SIMD/MIMD PMEs per chip. However, that can be extended to greater dimensionality or more PMEs per dimension as we will describe with respect to FIGURES 3, 9, 10, 15 and 16.

Our processing element includes a full I/O system including both data transfer and program interrupts. Our description of our preferred embodiment will be primarily described in terms of the preferred 4-dimensional I/O structures with 8 SIMD/MIMD PMEs per chip, which has special advantages now in our view. However, that can be extended to greater dimensionality or more PMEs per dimension as described in our parent application. In addition, for most applications we prefer and have made inventions in areas of greater dimensions with hypercube interconnections, preferably with the modified hypercube we describe. However, in some applications a 2-dimensional mesh interconnection of chips will be applicable to a task at hand. For instance, in certain military computers a 2 dimensional mesh will be suitable and cost effective.

This disclosure extends the concepts from the interprocessor communication to the external input/output facilities and describes the interfaces and modules required for control of the processing array. In summary three types of I/O, inter-processor, processors to/from external, and broadcast/control are described. Massively parallel processing systems require all these types of I/O bandwidth demands to be balanced with processor computing capability. Within the array these requirements will be satisfied by replicating a 16 bit (reduced) instruction set processor, augmented with very fast interrupt state swapping capability. That processor is referred to as the PME illustrating the preferred embodiment of our APAP. The characteristics of the PME are completely unique when compared with the processing elements on other massively parallel machines. It permits the processing, routing, storage and I/O to be completely distributed. This is not characteristic of any other design.

In a hypercube each PME can address as its neighbor, any PME whose address differs in any single bit position. In a ring, any PME can address as its neighbor the two PMEs whose addresses differ  $\pm 1$ . The modified hypercube of our preferred embodiment utilized for the APAP combines these approaches by building hypercubes out of rings. The intersection of rings is defined to be a node. Each node of our preferred system has its PME, memory and I/O, and other features of the node, formed in a semiconductor silicon low level CMOS DRAM chip. Nodes are constructed from multiple PMEs on each chip. Each PME exists in only one ring of nodes. PMEs within the node are connected by additional rings such that communications can be routed between rings within the node. This leads to the addressing structure where any PME can step messages toward the objective by addressing a PME in its own ring or an adjacent PME within the node. In essence a PME can address a PME whose address differs by 1 in one in the  $1n_2d$  bit field of its ring (where  $d$  is the number of PMEs in the ring) or the PME with the same address but existing in an adjacent dimension. The PME effectively appears to exist in  $n$  sets of rings, while in actuality it exists only in one real ring and one hidden ring totally contained within the chip. The dimensionality for the modified hypercube is defined to be the value  $n$  from the previous sentence.

We prefer to use a modified hypercube. This is elaborated in the part of this application describing the technology. Finally, PMEs within a ring are paired such that one moves data externally clockwise along a ring of nodes and the other moves data externally counterclockwise along the ring of nodes, thus dedicating a PME to an external port.

In our massively parallel machine, in our preferred embodiment, the interconnection and broadcast of data and instructions from one PME to another PME in the node and externally of the node to other nodes of a cluster or PMEs of a massively parallel processing environment are performed by a programmable router, allowing reconfiguration and virtual flexibility to the network operations. This important feature is fully distributed and embedded in the PME and allows for processor communication and data transfers among PMEs during operations of the system in SIMD and MIMD modes, as well as in the SIMD/MIMD and SMIMD modes of operation.

Within the rings each interconnection leg is a point-to-point connection. Each PME has a point-to-point connection with the two neighboring PMEs in its ring and with two neighboring PMEs in two adjacent rings. Three of these point-to-point connections are internal to the node, while the fourth point-to-point connection is to an adjacent node.

The massively parallel processing system uses the processing elements, with their local memory and interconnect topology to connect all processors to each other. Embedded within the PME is our fully

distributed I/O programmable router. Our system also provides an addition to the system which provides the ability to load and unload all the processing elements. With our zipper we provide a method for loading and unloading of the array of PEs and thus enable implementation of a fast I/O along an edge of the array's rings. To provide for external interface I/O any subset of the rings may be broken (un-zipped) across some dimension(s) with the resultant broken paths connected to the external interface. The co-pending application entitled "APAP I/O ZIPPER", filed concurrently herewith, USSN filed May 22, 1992,  
 5 describes our 'zipper' in additional detail. The 'zipper' can be applied to only the subset of links required to support the peak external I/O load, which in all configurations considered so far leads to its being applied only to one or two edges of the physical design.

10 The final type of I/O consists of data that must be broadcast to, or gathered from all PMEs, plus data which is too specialized to fit on the standard buses. Broadcast data includes commands, programs and data. Gathered data is primarily status and monitor functions while diagnostic and test functions are the specialized elements. Each node, in addition to the included set of PMEs, contains one Broadcast and Control Interface (BCI) section.

15 Consider PMEs interconnected in a modified 4 dimensional hypercube network. If each ring contains 16 PMEs, then the system will have 32,768 PMEs. The network diameter is 18 steps. Each PME contains in SAW the router and reconfiguration SAW to support a particular outgoing port. Thus, software routing provides the capability to reconfigure in the event of a faulty processing element or node. Inherent in a 4d, 25 Mhz network design with byte wide half duplex rings is the provision for 410 gigabytes per second peak  
 20 internal bandwidth.

The 4 dimensional hypercube leads to a particularly advantageous package. Eight of the PMEs (including data flow, memory and I/O paths and controls) are encompassed in a single chip. Thus, a node will be a single chip including pairs of elements along the rings. The nodes are configured together in an 8 X 8 array to make up a cluster. The fully populated machine is built up of an array of 8 X 8 clusters to  
 25 provide the maximum capacity of 32,768 PMEs.

Each PME is a powerful microcomputer having significant memory and I/O functions. There is multibyte data flow within a reduced instruction set (RISC) architecture. Each PME has 16 bit internal data flow and eight levels of program interrupts with the use of working and general registers to manage data flow. There is a circuit switched and store and forward mode for I/O transfer under PME software control. The SIMD  
 30 mode or MIMD mode is under PME software control. The PME can execute RISC instructions from either the BCI in a SIMD mode, or from its own main memory in MIMD mode. Specific RISC instruction code points can be reinterpreted to perform unique functions in the SIMD mode. Each PME can implement an extended Instruction Set Architecture and provide routings which perform macro level instructions such as extended precision fixed point arithmetic, floating point arithmetic, vector arithmetic, and the like. This  
 35 permits not only complex math to be handled but image processing activities for display of image data in multiple dimensions (2d and 3d images) and for multimedia applications. The system can select groups of PMEs for a function. PMEs assigned can allocate selected data and instructions for group processing. The operations can be externally monitored via the BCI. Each BCI has a primary control input, a secondary control input, and a status monitor output for the node. Within a node the 2n PMEs can be connected for a  
 40 binary hypercube communication network within the chip. Communication between PMEs is controlled by the bits in PME control registers under control of PME software. This permits the system to have a virtual routing capability. Each PME can step messages up or down its own ring or to its neighboring PME in either of two adjacent rings. Each interface between PMEs is a point-to-point connection. The I/O ports permit off-chip extensions of the internal ring to adjacent nodes of the system. The system is built up of  
 45 replications of a node to form a node array, a cluster, and other configurations.

To complement our system's SIMD, MIMD, SIMD/MIMD and SIM/MIMD functionality, our development we have provided unique operational modes. Among our SIMD/MIMD PME's unique modes are the new functional features referred to as the "store and forward / circuit switch" functions. These hardware functions  
 50 complemented with the on chip communication and programmable internal and external I/O routing provides the PME with very optimal data transferring capability. In preferred mode of operation the processor memory is generally the data sink for messages and data targeted at the PME in the store and forward mode. Messages and data not targeted for the PME are sent directly to the required output port when in circuit switched mode. The PME software performs the selected routing path while giving the PME a dynamically selectable store and forward / circuit switch functionality.

55 Among the advances we have provided is a fully distributed architecture for PMEs of a node. Each node has 2n processors, memory and I/O. Every PME will provide very flexible processing capability with 16 bit data flow, 64K bytes of local storage, store and forward/circuit switch logic, PME to PME communication, SIMD/MIMD switching capabilities, programmable routing, and dedicated floating point

EP 0 570 729 A2

assist logic. The organization of every PME and its communication paths with other PMEs within the same chip to minimize chip crossing delays. PME functions can be independently operated by the PME and integrated with functions in the node, a cluster, and larger arrays.

Our massively parallel system is made up of nodal building blocks of multiprocessor nodes, clusters of nodes, and arrays of PMEs already packaged in clusters. For control of these packaged systems we provide a system array director which with the hardware controllers performs the overall Processing Memory Element (PME) Array Controller functions in the massively parallel processing environment. The Director comprises of three functional areas, the Application Interface, the Cluster Synchronizer, and normally a Cluster Controller. The Array Director will have the overall control of the PME array, using the broadcast bus and our zipper connection to steer data and commands to all of the PMEs. The Array Director functions as a software system interacting with the hardware to perform the role as the shell of the APAP operating system.

The interconnection for our PMEs for a massively parallel array computer SIMD/MIMD processing memory element (PME) interconnection provides the processor to processor connection in the massively parallel processing environment. Each PME utilizes our fully distributed interprocessor communication hardware from the on-chip PME to PME connection, to the off-chip I/O facilities which support the chip-to-chip interconnection. Our modified topology limits our cluster to cluster wiring while supporting the advantages of hypercube connections.

The concepts which we employ for a PME node are related to the VLSI packaging techniques used for the Advanced Parallel Array Processor (APAP) computer system disclosed here, which packaging features of our invention provide enhancements to the manufacturing ability of the APAP system. These techniques are unique in the area of massively parallel processor machines and will enable the machine to be packaged and configured in optimal subsets that can be built and tested.

The packaging techniques take advantage of the eight PMEs packaged in a single chip and arranged in a N-dimensional modified hypercube configuration. This chip level package or node of the array is the smallest building block in the APAP design. These nodes are then packaged in an 8 x 8 array where the +-X and the +-Y makes rings within the array or cluster and the +-W, and +-Z are brought out to the neighboring clusters. A grouping of clusters make up an array. The intended applications for APAP computers depend upon the particular configuration and host. Large systems attached to mainframes with effective vectorized floating point processors might address special vectorizable problems - such as weather prediction, wind tunnel simulation, turbulent fluid modeling and finite element modeling. Where these problems involve sparse matrices, significant work must be done to prepare the data for vectorized arithmetic and likewise to store results. That workload would be off loaded to the APAP. In intermediate size systems, the APAP might be dedicated to performing the graphics operations associated with visualization, or with some preprocessing operation on incoming data (i.e., performing optimum assignment problems in military sensor fusion applications). Small systems attached to workstations or PCs might serve as programmer development stations or might emulate a vectorized floating point processor attachment or a 3d graphics processor.

#### 40 BRIEF DESCRIPTION OF THE DRAWINGS

- FIG. 1 shows a parallel processor processing element like those which would utilize old technology.
- FIG. 2 shows a massively parallel processor building block in accordance with our invention, representing our new chip design.
- 45 FIG. 3 illustrates on the right side the preferred chip physical cluster layout for our preferred embodiment of a chip single node fine grained parallel processor. There each chip is a scalable parallel processor chip providing 5 MIPS performance with CMOS DRAM memory and logic permitting air cooled implementation of massive concurrent systems. On the left side of Figure 3, there is illustrated the replaced technology.
- 50 FIG. 4 shows a computer processor functional block diagram in accordance with the invention.
- FIG. 5 shows a typical Advanced Parallel Array Processor computer system configuration.
- FIG. 6 shows a system overview of our fine-grained parallel processor technology in accordance with our invention, illustrating system build up using replication of the PME element which permits systems to be developed with 40 to 193,840 MIPS performance.
- 55 FIG. 7 illustrates the hardware for the processing element (PME) data flow and local memory in accordance with our invention, while
- FIG. 8 illustrates PME data flow where a processor memory element is configured as a hardwired general purpose computer that provides about 5 MIPS fixed point processing or 4 Mflops via

- programmed control floating point operations.
- FIG. 9 shows the PME to PME connection (binary hypercube) and data paths that can be taken in accordance with our invention, while
- FIG. 10 illustrates node interconnections for the chip or node which has 8 PMEs, each of which manages a single external port and permits distribution of the network control function and eliminates a functional hardware port bottleneck.
- FIG. 11 is a block diagram of a scalable parallel processor chip where each PME is a 16 bit wide processor with 32 K words of local memory and there is I/O porting for a broadcast port which provides a controller-to-all interface while external ports are bi-directional point-to-point interfaces permitting ring torus connections within the chip and externally.
- FIG. 12 shows an array director in the preferred embodiment.
- FIG. 13 in part (a) illustrates the system bus to or from a cluster array coupling enabling loading or unloading of the array by connecting the edges of clusters to the system bus (see FIGURE 14). In FIGURE 13 in part (b) there is the bus to/from the processing element portion. FIGURE 13 illustrates how multiple system buses can be supported with multiple clusters. Each cluster can support 50 to 57 Mbytes/s bandwidth.
- FIG. 14 shows a "zipper" connection for fast I/O connection.
- FIG. 15 shows an 8 degree hypercube connection illustrating a packaging technique in accordance with our invention applicable to an 8 degree hypercube.
- FIG. 16 shows two independent node connections in the hypercube.
- FIG. 17 shows the Bitonic Sort algorithm as an example to illustrate the advantages of the defined SIMD/MIMD processor system.
- FIG. 18 illustrates a system block diagram for a host attached large system with one application processor interface illustrated. This illustration may also be viewed with the understanding that our invention may be employed in stand alone systems which use multiple application processor interfaces. Such interfaces in a FIGURE 18 configuration will support DASD/Graphics on all or many clusters. Workstation accelerators can eliminate the host, application processor interface (API) and cluster synchronizer (CS) illustrated by emulation. The CS is not required in all instances.
- FIG. 19 illustrates the software development environment for our system. Programs can be prepared by and executed from the host application processor. Both program and machine debug is supported by the workstation based console illustrated here and in FIGURE 22. Both of these services will support applications operating on a real or a simulated MMP, enabling applications to be developed at a workstation level as well as on a supercomputer formed of the APAP MMP. The common software environment enhances programmability and distributed usage.
- FIG. 20 illustrates the programming levels which are permitted by the new systems. As different users require more or less detailed knowledge, the software system is developed to support this variation. At the highest level the user does not need to know the architecture is indeed an MMP. The system can be used with existing language systems for partitioning of programs, such as parallel Fortran.
- FIG. 21 illustrates the parallel Fortran compiler system for the MMP provided by the APAP configurations described. A sequential to parallel compiler system uses a combination of existing compiler capability with new data allocation functions and enables use of a partitioning program like FortranD.
- FIG. 22 illustrates the workstation application of the APAP, where the APAP becomes a workstation accelerator. Note that the unit has the same physical size as a RISC/6000 Model 530, but this model now contains an MMP which is attached to the workstation via a bus extension module illustrated.
- FIG. 23 illustrates an application for an APAP MMP module for an AWACS military or commercial application. This is a way of handling efficiently the classical distributed sensor fusion problem shown in FIGURE 23, where the observation to track matching is classically done with well know algorithms like nearest neighbor, 2 dimensional linear assignment (Hunges...), probabilistic data association or multiple hypothesis testing, but these can now be done in an improved manner as illustrated by FIGURES 24 and 25.
- FIG. 24 illustrates how the system provides the ability to handle n-dimensional assignment problems in real time.
- FIG. 25 illustrates processing flow for an n-dimensional assignment problem utilizing an APAP.

EP 0 570 729 A2

FIG. 26 illustrates the expansion unit provided by the system enclosure described showing how a unit can provide 424 MflopS or 5120 MIPs using only 8 to 10 extended SEM-E modules, providing the performance comparable to that of specialized signal processor module in only .8 cubic feet. This system can become a SIMD massive machine with 1024 parallel processors performing two billion operations per second (GOPS) and can grow by adding 1024 additional processors and 32MB additional storage.

FIG. 27 illustrates the APAP packaging for a supercomputer. Here is a large system of comparable performance but much smaller footprint than other systems. It can be built by replicating the APAP cluster within an enclosure like those used for smaller machines.

We have provided, as part of the description, Tables illustrating the hardwired instructions for a PME, in which Table 1 illustrates Fixed-point arithmetic instructions; Table 2 illustrates storage to storage instructions; Table 3 illustrates logical instructions; Table 4 illustrates shift instructions; Table 5 illustrates branch instructions; Table 6 illustrates the status switching instructions; and Table 7 illustrates the input/output instructions.

(Note: For convenience of illustration in the formal patent drawings, FIGURES may be separated in parts and as a convention we place the top of the FIGURE as the first sheet, with subsequent sheets proceeding down and across when viewing the FIGURE, in the event that multiple sheets are used.)

Our detailed description follows with parts explaining the preferred embodiments of our invention provided by way of example.

#### DETAILED DESCRIPTION OF THE INVENTION

Turning now to our invention in greater detail, it will be seen from FIGURE 1, which illustrates the existing technology level, illustrated by the transputer T000 chip, and representing similar chips for such machines as the illustrated by the Touchstone Delta (1980), N Cube ('386), and others. When FIGURE 1 is compared with the developments here, it will be seen that not only can systems like the prior systems be substantially improved by employing our invention, but also new powerful systems can be created, as we will describe. FIGURE 1's conventional modern microprocessor technology consumes pins and memory. Bandwidth is limited and inter-chip communication drags the system down.

The new technology leapfrog represented by FIGURE 2 merges processors, memory, I/O into multiple PMEs (eight or more 16 bit processors each of which has no memory access delays and uses all the pins for networking) formed on a single low power CMOS DRAM chip. The system can make use of ideas of our prior referenced disclosures as well as invention separately described in the applications filed concurrently herewith and applicable to the system we describe here. Thus, for this purpose they are incorporated herein by reference. Our concepts of grouping, autonomy, transparency, zipper interaction, asynchronous SIMD, SIMIMD or SIMD/MIMD, can all be employed with the new technology, even though to lesser advantage they can be employed in the systems of the prior technology and in combination with our own prior multiple picket processor.

Our picket system can employ the present processor. Our basic concept is that we have now provided a replicable brick, a new basic building block for systems with our new memory processor, a memory unit having embedded processors, router and I/O. This basic building block is scalable. The basic system which we have implemented employs a 4 Meg. CMOS DRAM. It is expandable to be used in larger memory configurations, with 16Mbit DRAMS, and 64Mbit chips by expansion. Each processor is a gate array. With denser deposition, many more processors, at higher clock speeds, can be placed on the same chip, and using gates and additional memory will expand the performance of each PME. Scaling a single part type provides a system framework and architecture which can have a performance well into the PETAOP range.

FIGURE 2 illustrates the memory processor which we call the PME or processor memory element in accordance with our preferred embodiment. The processor has eight or more processors. In the pictured embodiment there are eight. The chip can be expanded (horizontally) to add more processors. The chip can, as preferred, retain the logic and expand the DRAM memory with additional cells linearly (vertically). Pictured are 16 - 32k by 9 bit sections of DRAM memory surrounding a field of CMOS gate array gates which implement 8 replications of a 16 bit wide data flow processors.

Using IBM CMOS low power sub-micron IBM CMOS deposition on silicon technology, it uses selected silicon with trench to provide significant storage on a small chip surface. Our memory and multiple processors organized interconnect is made with IBM's advanced art of making semiconductor chips. However, it will be recognized that the little chip we describe has about 4 Meg. memory. It is designed so that as 16 Meg. memory technology becomes stable, when improved yields and methods of accommodating defects are certain, our little chip can migrate to larger memory sizes each 9 bits wide without changing

the logic. Advances in photo and X-ray lithography keep pushing minimum feature size to well below .5 microns. Our design envisions more progress. These advances will permit placement of very large amounts of memory with processing on a single silicon chip.

Our device is a 4 MEG CMOS DRAM believed to be the first general memory chip with extensive room for logic. 16 replications of a 32k by 9-bit DRAM macro make up the memory array. The DRAM has 120K cells it allocates with significant surface area for application logic on the chip, with triple level metal wiring. The processor logic cells are preferably gate array cells. The 35 ns or less DRAM access time matches the processor cycle time. This CMOS implementation provides logic density for a very effective PE (picket) and does so while dissipating 1.3 watts for the logic. The separate memory section of the chip, each 32K by 9 bits, (with expansion not changing logic) surrounds the field of CMOS gate array gates representing 120K cells, and having the logic described in other figures. Memory is barriered and with a separated power source dissipates .9 watts. In providing the combining of significant amounts of logic on the same silicon substrate with significant amounts of memory problems involved with the electrical noise incompatibility of logic and DRAM have been overcome. Logic tends to be very noisy while memory needs relative quiet to sense the millivolt size signals that result from reading the cells of DRAM. We prefer to provide trenched triple metal layer silicon deposition, with separate barriered portions of the memory chip devoted to memory and to processor logic with voltage and ground isolation, and separate power distribution and barriers, to achieve compatibility between logic and DRAM.

## 20 APAP System Overview of Preferred Embodiments

This description introduces the new technology in the following order:

1. Technology
2. Chip H/W description
- 25 3. Networking and system build up
4. Software
5. Applications

The initial sections of the detailed description describe how 4-Meg DRAM low power CMOS chips are made to include 8 processors on and as part of the manufactured PME DRAM chips each supporting:

- 30 1. 16 bit, 5 MIP dataflows,
2. independent instruction stream and interrupt processing and
3. 8 bit (plus parity and controls) wide external port and interconnection to 3 other on chip processors.

Our invention provides multiple functions which are integrated into a single chip design. The chip will provide PME functions which are powerful and flexible and sufficiently so such that a chip having scalability will be effective at processing, routing, storage and three classes of I/O. This chip has integrated memory and control logic within the single chip to make the PME, and this combination is replicated within the chip. A processor system is built from replications of the single chip.

The approach partitions the low power CMOS DRAM. It will be formed as multiple word length (16) bit by 32K sections, associating one section with a processor. (We use the term PME to refer to a single processor, memory and I/O capable system unit.) This partitioning leads to each DRAM chip being an 8 way 'cube connected' MIMD parallel processor with 8 byte wide independent interconnection ports. (See FIGURE 6 for an illustration of a replication of fine-grained parallel technology. Illustrating replication and the ring torus possibilities.)

The software description addresses several distinct program types. At the lowest level, processes interface the user's program (or services called by the application) to the detailed hardware H/W needs. This level includes the tasks required to manage the I/O and interprocessor synchronization and is what might be called a microprogram for the MPP. An intermediate level of services provide for both mapping applications (developed with vector or matrix operations) to the MPP, and also control, synchronization, startup, diagnostic functions. At the host level, high order languages are supported by library functions that support vectorized programs with either simple automatic data allocation to the MPP or user tuned data allocation. The multi-level software S/W approach permits applications to exploit different degrees of control and optimization within a single program. Thus, a user can code application programs without understanding the architecture detail while an optimizer might tune at the microcode level only the small high usage kernels of a program.

Sections of our description that describe 1024 element 6 GIPS units and a 32,768 element 164 GIPS unit illustrate the range of possible systems. However, those are not the limits; both smaller and larger units are feasible. These particular sizes have been selected as examples because the small unit is suitable to microprocessors (accelerators), personal computers, workstation and military applications (using of coarse



different packaging techniques), while the larger unit is illustrative of a mainframe application as a module or complete supercomputer system. A software description will provide examples of other challenging work that might be effectively programmed on each of the illustrative systems.

#### 9 PME DRAM CMOS - A BASE FOR A MULTIPROCESSOR PME

FIGURE 2 illustrates our technology improvement at the chip technology level. This extendable computer organization is very cost and performance efficient over the wide range of system sizes because it uses only one chip type. Combining the memory and processing on one chip eliminates the pins  
10 dedicated to the memory bus and their associated reliability and performance penalties. Replication of our design within the chip makes it economically feasible to consider custom logic designs for processor subsections. Replication of the chip within the system leads to large scale manufacturing economies. Finally, CMOS technology requires low power per MIP, which in turn minimizes power supply and cooling needs. The chip architecture can be programmed for multiple word lengths enabling operations to be  
15 performed that would otherwise require much larger length processors. In combination these attributes permit the extensive range of system performance.

Our new technology can be compared with a possible extension of the old technology it overlaps. It is apparent that the advantages of smaller features have been used by processor designers to construct more complex chips and by memory designers to provide greater replication of the simple element. If the trend  
20 continues one could expect memories to get four times as large while processors might exploit density to:

1. include multiple execute units with instruction routers,
2. increase cache sizes and associative capability and/or
3. increase instruction look ahead and advance computation capability.

However, these approaches to the old technology illustrated by FIGURE 1 all tend to dead end.  
25 Duplicating processors leads to linearly increasing pin requirements but pins per chip is fixed. Better caching can only exploit the application's data reuse pattern. Beyond that, memory bandwidth becomes the limit. Application data dependencies and branching limit the potential advantage of look ahead schemes. Additionally, it is not apparent that MPP applications with fine-grained parallelism need 1, 4, or 16 Megaword memories per processing unit. Attempting to share such large memories between multiple  
30 processors results in severe memory bandwidth limitations.

Our new approach is not dead ended. We combine both significant memory and I/O and processor into a single chip, as illustrated by the FIGURE 2 and subsequent illustration and description. It reduces part number requirements and eliminates the delays associated with chip crossing. More importantly, this permits all the chip's I/O pins to be dedicated to interprocessor communication and thus, maximizes  
35 network bandwidth.

To implement our preferred embodiment illustrated in FIGURE 2 we use a process that is available now, using IBM low power CMOS technology. Our illustrated embodiment can be made with CMOS DRAM density, in CMOS and can be implemented in denser CMOS. Our illustrated embodiment of 32K memory cells for each of 8 PHEs on a chip can be increased as CMOS becomes denser. In our embodiment we  
40 utilize the real estate and process technology for a 4 MEG CMOS DRAM, and expand this with processor replication associated with 32K memory on the chip itself. The chip, it will be seen has processor, memory, and I/O in each of the chip packages of the cluster shown in FIGURE 3. Within each package is a memory with embedded processor element, router, and I/O, all contained in a 4 MEG CMOS DRAM believed to be the first general memory chip with extensive room for logic. It uses selected silicon with trench to provide  
45 significant storage on a small chip surface. Each processor chip of our design alternatively can be made with 16 replications of a 32K by 9 bit DRAM macro (35/80 ns) using .07 micron CMOS logic to make up the memory array. The device is unique in that it allocates surface area for 120 K cells of application logic on the chip, supported by the capability of triple level metal wiring. The multiple cards of the old technology is shown crossed out on the left side of FIGURE 3.

Our basic replicable element brick technology is an answer to the old technology. If one considered the "Xed" technology on the left of FIGURE 3, one would see too many chips, too many cards, and waste. For example, today's proposed teraflop machines that others offer would have literally a million or more chips in them. With today's other technology only a few percent of these chips, at best, are truly operations producers. The rest are "overhead" (typically memory, network interface, etc.).  
50

It will become evident that it is not feasible to package such chips, in such a large number, in anything that must operate in a constrained environment of physical size. (How many could you fit in a small area of a cockpit?) Furthermore, such proposed teraflop machines or others, already huge, must scale up 1000x times to reach the petaop range. We have a solution which dramatically decreases the percent of non-  
55

operations producing chips. We provide increased bandwidth. We provide this within a reasonable network dimensionality. With such a brick technology, where memory becomes the operator, and networks are used for passing controls, where operations producing chips are dramatically increased. In addition, the upgrade dramatically reduces the number of different types of chips. Our system is designed for scale-up, without a requirement for specialized packaging, cooling, power, or environmental constraints.

With our brick technology, utilizing instead of separate processors, memory units with built in processors and network capability, the configuration shown in FIGURE 3, representing a card, with chips which are pin compatible with current 4Mbit DRAM cards at the connector level. Such a single card could hold, with a design point of a basic 40 mips per chip performance level, 32 chips, or 1280 mips. Four such cards would provide 5 gips. The workstation configuration which is illustrated would preferably have such a PE memory array, a cluster controller, and an IBM RISC System/6000 which has sufficient performance to run and monitor execution of an array processor application developed at the workstation.

A very gate efficient processor can be used in the processor portion. Such designs for processors have been employed, but never within memory. Indeed, in addition, we have provided the ability to mix MIMD and SIMD basic operation provisions. Our chip provides a "broadcast bus" which provides an alternate path into each CPU's instruction buffer. Our cluster controller issues commands to each of the PEs in the PMEs, and these can be stored in the PME to control their operation in one mode or another. Each PME does not have to store an entire program, but can store only those portions applicable to a given task at various times during processing of an application.

Given the basic device one can elect to develop a single processor memory combination. Alternatively, by using a more simple processor and a subset of the memory macros one can design for either 2, 4, 8 or 16 replications of the basic processing element (PME). The PME can be made simpler either by adjusting the dataflow bandwidth or by substituting processor cycles for functional accelerators. For most embodiments we prefer to make 8 replications of the basic processing element we describe.

Our application studies have indicated that for now the most favorable answer is 8 replications of a 16 bit wide data flow and 32K word memory. We conclude this because:

1. 16 bit words permit single cycle fetch of instructions and addresses.
2. 8 PMEs each with an external port permits 4 dimensional torus interconnections. Using 4 or 8 PMEs on each ring leads to modules suitable for the range of targeted system performances.
3. 8 external ports requires about 60% of the chip pins, providing sufficient remainder for power, ground and common control signals.
4. 8 Processors implemented in a 64 KByte Main Store
  - a. allows for a register based architecture rather than a memory mapped architecture, and it
  - b. forces some desirable but not required accelerators to be implemented by multiple processor cycles.

This last attribute is important because it permits use of the developing logic density increase. Our new accelerators (ex. floating point arithmetic unit per PME) are added as chip hardware without affecting system design, pins and cables or application code.

The resultant chip layout and size (14.59 x 14.83 mm) is shown in FIGURE 2, and FIGURE 3 shows a cluster of such chips, which can be packaged in systems like those shown in later FIGURES for stand alone units, workstations which slide next to a workstation host with a connection bus, in AWACs applications, and in supercomputers. This chip technology provides a number of system level advantages. It permits development of the scalable MPP by basic replication of a single part type. The two DRAM macros per processor provide sufficient storage for both data and program. An SRAM of equivalent size might consume more than 10 times more power. This advantage permits MIMD machine models rather than the more limited SIMD models characteristic of machines with single chip processor/memory designs. The 35 ns or less DRAM access time matches the expected processor cycle time. CMOS logic provides the logic density for a very effective PME and does so while dissipating only 1.3 watts. (Total chip power is 1.3 + .9 (memory) = 2.2 w.) Those features in turn permit using the chip in MIL applications requiring conduction cooling. (Air cooling in non-MIL applications is significantly easier.) However, the air cooled embodiment can be used for workstation and other environments. A stand-alone processor might be configured with an 80 amp - 5 volt power supply.

Advanced Parallel Array Processor (APAP) building blocks are shown in FIGURE 4 and in FIGURE 5. FIGURE 4 illustrates the functional block diagram of the Advanced Parallel Array Processor. Multiple application interfaces 150, 160, 170, 180 exist for the application processor 100 or processors 110, 120, 130. FIGURE 5 illustrates the basic building blocks that can be configured into different system block diagrams. The APAP, in a maximum configuration, can incorporate 32,768 identical PMEs. The processor consists of the PME Array 280, 290, 300, 310, an Array Director 250 and an Application Processor Interface

260 for the application processor 200 or processors 210, 220, 230. The Array Director 250 consists of three functional units: Application Processor Interface 260, cluster Synchronizer 270 and cluster Controller 270. An Array Director can perform the functions of the array controller of our prior linear picket system for SIMD operations with MIMD capability. The cluster controller 270, along with a set of 64 Array clusters 280, 290, 300, 310, (i.e. cluster of 512 PMEs), is the basic building block of the APAP computer system. The elements of the Array Director 250 permit configuring systems with a wide range of cluster replications. This modularity based upon strict replication of both processing and control elements is unique to this massively parallel computer system. In addition, the Application Processor Interface 260 supports the Test/Debug device 240 which will accomplish important design, debug, and monitoring functions.

Controllers are assembled with a well-defined interface, e.g. IBM's Microchannel, used in other systems today, including controllers with i860 processors. Field programmable gate arrays add functions to the controller which can be changed to meet a particular configuration's requirements (how many PMEs there are, their couplings, etc.)

The PME arrays 280, 290, 300, 310 contain the functions needed to operate as either SIMD or MIMD devices. They also contain functions that permit the complete set of PMEs to be divided into 1 to 256 distinct subsets. When divided into subsets the Array Director 250 interleaves between subsets. The sequence of the interleave process and the amount of control exercised over each subset is program controlled. This capability to operate distinct subsets of the array in one mode, i.e., MIMD with differing programs, while other sets operate in tightly synchronized SIMD mode under Array Director control, represents an advance in the art. Several examples presented later illustrate the advantages of the concept.

#### Array Architecture

The set of nodes forming the Array is connected as a n-dimensional modified hypercube. In that interconnection scheme, each node has direct connections to  $2n$  other nodes. Those connections can be either simplex, half-duplex or full-duplex type paths. In any dimension greater than 3d, the modified hypercube is a new concept in interconnection techniques (The modified hypercube in the 2d case generates a torus, and in the 3d case an orthogonally connected lattice with edge surfaces wrapped to opposing surface.)

To describe the interconnection scheme for greater than 3d cases requires an inductive description. A set of  $m_1$  nodes can be interconnected as a ring. (The ring could be 'simply connected', 'braided', 'cross connected', 'fully connected', etc. Although additional node ports are needed for greater than simple rings, that added complexity does not affect the modified hypercube structure.) The  $m_1$  rings can then be linked together by connecting each equivalent node in the  $m_2$  set of rings. The result at this point is a torus. To construct a  $i+1$ d modified hypercube from an  $i$ d modified hypercube,  $m_{i+1}$  sets of  $i$ d modified hypercubes and interconnect all of the equivalent  $m_i$  level nodes into rings.

This process is illustrated for the 4d modified hypercube, using  $m_i = 8$  for  $i = 1..4$  by the illustration in FIGURE 6. Compare our description under node Topology and also FIGURES 8, 9, 10, 15 and 18.

FIGURE 8 illustrates the fine-grained parallel technology path from the single processor element 300, made up of 32K 16-bit words with a 16-bit processor to the Network node 310 of eight processors 312 and their associated memory 311 with their fully distributed I/O routers 313 and Signal I/O ports 314, 315, on through groups of nodes labeled clusters 320 and into the cluster configuration 360 and to the various applications 330, 340, 350, 370. The 2d level structure is the cluster 320, and 64 clusters are integrated to form the 4d modified hypercube of 32,768 Processing Elements 360.

#### Processing Array Element (PME) Preferred Embodiment

As illustrated by FIGURE 2 and FIGURE 11 the preferred APAP has a basic building block of a one chip node. Each node contains 8 identical processor memory elements (PMEs) and one broadcast and control interface (BCI). While some of our inventions may be implemented when all functions are not on the same chip, it is important from a performance and cost reduction standpoint to provide the chip as a one chip node with the 8 processor memory elements using the advanced technology which we have described and can be implemented today.

The preferred implementation of a PME has a 64 KByte main store, 16 16-bit general registers on each of 8 program interrupt levels, a full function arithmetic/logic unit (ALU) with working registers, a status register, and four programmable bi-directional I/O ports. In addition the preferred implementation provides a SIMD mode broadcast interface via the broadcast and control interface (BCI) which allows an external controller (see our original parent application and the description of our currently preferred embodiment for

a nodal array and system with clusters) to drive PME operation decode, memory address, and ALU data inputs. This chip can perform the functions of a microcomputer allowing multiple parallel operations to be performed within it, and it can be coupled to other chips within a system of multiple nodes, whether by an interconnection network, a mesh or hypercube network, or our preferred and advanced scalable embodiment.

The PMEs are interconnected in a series of rings or tori in our preferred scalable embodiment. In some applications the nodes could be interconnected in a mesh. In our preferred embodiment each node contains two PMEs in each of four tori. The tori are denoted W,X,Y, and Z (see FIGURE 8). FIGURE 11 depicts the interconnection of PMEs within a node. The two PMEs in each torus are designated by their external I/O port (+W, -W, +X, -X, +Y, -Y, +Z, -Z). Within the node, there are also two rings which interconnect the 4 +n and 4 -n PMEs. These internal rings provide the path for messages to move between the external tori.

Since the APAP can be in our preferred embodiment a four dimensional orthogonal array, the internal rings allow messages to move throughout the array in all dimensions.

The PMEs are self-contained stored program microcomputers comprising a main store, local store, operation decode, arithmetic/logic unit (ALU), working registers and Input/Output I/O ports. The PMEs have the capability of fetching and executing stored instructions from their own main store in MIMD operation or to fetch and execute commands via the BCI interface in SIMD mode. This interface permits intercommunication among the controller, the PME, and other PMEs in a system made up of multiple chips.

The BCI is the node's interface to the external array controller element and to an array director. The BCI provides common node functions such as timers and clocks. The BCI provides broadcast function masking for each nodal PME and provides the physical interface and buffering for the broadcast-bus-to-PME data transfers, and also provides the nodal interface to system status and monitoring and debug elements.

Each PME contains separate interrupt levels to support each of its point-to-point interfaces and the broadcast interface. Data is input to the PME main store or output from PME main store under Direct Memory Access (DMA) control. An "input transfer complete" interrupt is available for each of the interfaces to signal the PME software that data is present. Status information is available for the software to determine the completion of data output operations.

Each PME has a "circuit switched mode" of I/O in which one of its four input ports can be switched directly to ones of its four output ports, without having the data enter the PME main store. Selection of the source and destination of the "circuit switch" is under control of the software executing on the PME. The other three input ports continue to have access to PME main store functions, while the fourth input is switched to an output port.

An additional type of I/O has data that must be broadcast to, or gathered from all PMEs, plus data which is too specialized to fit on the standard buses. Broadcast data can include SIMD commands, MIMD programs, and SIMD data. Gathered data is primarily status and monitor functions. Diagnostic and test functions are the specialized data elements. Each node, in addition to the included set of PMEs, contains one BCI. During operations the BCI section monitors the broadcast interface and stores/collects broadcast data to/from the addressed PME(s). A combination of enabling masks and addressing tags are used by the BCI to determine what broadcast information is intended for which PMEs.

Each PME is capable of operating in SIMD or in MIMD mode in our preferred embodiment. In SIMD mode, each instruction is fed into the PME from the broadcast bus via the BCI. The BCI buffers each broadcast data word until all of its selected nodal PMEs have used it. This synchronization provides accommodation of the data timing dependencies associated with the execution of SIMD commands and allows asynchronous operations to be performed by a PME. In MIMD mode, each PME executes its own program from its own main store. The PMEs are initialized to the SIMD mode. For MIMD operations, the external controller normally broadcasts the program to each of the PMEs while they are in SIMD mode, and then commands the PMEs to switch to MIMD mode and begin executing. Masking/tagging the broadcast information allows different sets of PMEs to contain different MIMD programs, and/or selected sets of PMEs to operate in MIMD mode while other sets of PMEs execute in SIMD mode. In various software clusters or partitions these separate functions can operate independently of the actions in other clusters or partitions. The operation of the Instruction Set Architecture (ISA) of the PME will vary slightly depending on whether the PME is in the SIMD or MIMD mode. Most ISA instructions operate identically regardless of mode. However, since the PME in SIMD mode does not perform branching or other control functions some code points dedicated to those MIMD instructions are reinterpreted in SIMD mode to allow the PME to perform special operations such as searching main memory for a match to a broadcast data value or switching to MIMD mode. This further extends system flexibility of an array.

PME Architecture

Basically, our preferred architecture comprises a PME which has a 16 bit wide data flow, 32K of 16 bit memory, specialized I/O ports and I/O switching paths, plus the necessary control logic to permit each PME to fetch, decode and execute the 16 bit instruction set provided by our instruction set architecture (ISA). The preferred PME performs the functions of a virtual router, and thus performs both the processing functions and data router functions. The memory organization allows by cross addressing of memory between PMEs access to a large random access memory, and direct memory for the PME. The individual PME memory can be all local, or divided into local and shared global areas programmatically. Specialized controls and capabilities which we describe permit rapid task switching and retention of program state information at each of the PMEs interrupt execution levels. Although some of the capabilities we provide have existed in other processors, their application for management of interprocessor I/O is unique in massively parallel machines. An example is the integrate of the message router function into the PME itself. This eliminates specialized router chips or development of specialized VLSI routers. We also recognize that in some instances one could distribute the functions we provide on a single chip onto several chips interconnected by metalization layers or otherwise and accomplish improvements to massively parallel machines. Further, as our architecture is scalable from a single node to massively parallel supercomputer level machines, it is possible to utilize some of our concepts at different levels. As we will illustrate for example our PME data flow is very powerful, and yet operates to make the scalable design effective.

The PME processing memory element develops for each of the multiple PMEs of a node, a fully distributed architecture. Every PME will be comprised of processing capability with 16 bit data flow, 64K bytes of local storage, store and forward/circuit switch logic, PME to PME communication, SIMD/MIMD switching capabilities, programmable routing, and dedicated floating point assist logic. These functions can be independently operated by the PME and integrated with other PMEs within the same chip to minimize chip crossing delays. Referring to FIGURES 7 and 8 we illustrate the PME dataflow. The PME consists of 16 bit wide dataflow 425, 435, 445, 455, 465, 32K by 16 bit memory 420, specialized I/O ports 400, 410, 480, 490 and I/O switching paths 425, plus the necessary control logic to permit the PME to fetch, decode and execute a 16 bit reduced instruction set 430, 440, 450, 480. The special logic also permits the PME to perform as both the processing unit 460 and data router. Specialized controls 405, 406, 407, 408 and capabilities are incorporated to permit rapid task switching and retention of program state information at each of the PMEs' interrupt execution levels. Such capabilities have been included in other processors; however, their application specifically for management of interprocessor I/O is unique in massively parallel machines. Specifically, it permits the integration of the router function into the PME without requiring specialized chips or VLSI development macros.

16 bit internal data flow and control

The major parts of the internal data flow of the processing element are shown in FIGURE 7. FIGURE 7 illustrates the internal data flow of the processing element. This processing element has a full 16 bit internal data flow 425, 435, 445, 455, 465. The important paths of the internal data flows use 12 nanosecond hard registers such as the OP register 450, M register 440, WR register 470, and the program counter PC register 430. These registers feed the fully distributed ALU 460 and I/O router registers and logic 405, 406, 407, 408 for all operations. With current VLSI technology, the processor can execute memory operations and instruction steps at 25 Mhz, and it can build the important elements, OP register 450, M register 440, WR register 470, and the PC register 430 with 12 nanosecond hard registers. Other required registers are mapped to memory locations.

As seen in FIGURE 8 the internal data flow of the PME has its 32K by 16 bit main store in the form of two DRAM macros. The remainder of the data flow consists of CMOS gate array macros. All of the memory can be formed with the logic with low power CMOS DRAM deposition techniques to form an very large scaled integrated PME chip node. The PME is replicated 8 times in the preferred embodiment of the node chip. The PME data flow consists of a 16 word by 16 bit general register stack, a multi-function arithmetic/logic unit (ALU) working registers to buffer memory addresses, memory output registers, ALU output registers, operation/command, I/O output registers, and multiplexors to select inputs to the ALU and registers. Current CMOS VLSI technology for 4MByte DRAM memory with our logic permits a PME to execute instruction steps at 25Mhz. We are providing the OP register, the M register, the WR register and the general register stack with 12 nanosecond hard registers. Other required registers are mapped to memory locations within a PME.

The PME data flow is designed as a 16 bit integer arithmetic processor. Special multiplexor paths have been added to optimize subroutine emulation of  $n \times 16$  bit floating point operations ( $n \geq 1$ ). The 16 bit data flow permits effective emulation of floating point operations. Specific paths within the data flow have been included to permit floating point operations in as little as 10 cycles. The ISA includes special code point to permit subroutines for extended (longer than 16-bit) operand operations. The subsequent floating point performance is approximately one twentieth the fixed floating point performance. This performance is adequate to eliminate the need for special floating point chips augmenting the PME as is characteristic of other massively parallel machines. Some other processors do include the special floating point processors on the same chip as a single processor (See FIGURE 1). We can enable special floating point hardware processors on the same chip with our PMEs but we would now need additional cells than is required for the preferred embodiment. For floating point operations, see also the concurrently filed FLOATING POINT application referenced above for improvements to the IEEE standard.

The approach developed is well poised to take advantage of the normal increases in VLSI technology performance. As circuit size shrinks and greater packaging density becomes possible then data flow elements like base and index registers, currently mapped to memory could be moved to hardware. Likewise, floating point sub-steps are accelerated with additional hardware which we will prefer for the developing CMOS DRAM technology as reliable higher density levels are achieved. Very importantly, this hardware alternative does not affect any software.

The PME is initialized to SIMD mode with interrupts disabled. Commands are fed into the PME operation decode buffer from the BCI. Each time an instruction operation completes, the PME requests a new command from the BCI. In a similar manner, immediate data is requested from the BCI at the appropriate point in the instruction execution cycle. Most instructions of the ISA operate identically whether the PME is in SIMD mode or in MIMD mode, with the exception of that SIMD instructions and immediate data are taken from the BCI; in MIMD mode the PME maintains a program counter (PC) and uses that as the address within its own memory to fetch a 16 bit instruction. Instructions such as "Branch" which explicitly address the program counter have no meaning in SIMD mode and some of those code points are reinterpreted to perform special SIMD functions as comparing immediate data against an area of main store.

The PME instruction decode logic permits either SIMD/MIMD operational modes, and PMEs can transition between modes dynamically. In SIMD mode the PME receives decoded instruction information and executes that data in the next clock cycle. In MIMD mode the PME maintains a program counter PC address and uses that as the address within its own memory to fetch a 16 bit instruction. Instruction decode and execution proceeds as in most any other RISC type machine. A PME in SIMD mode enters MIMD mode when given the information that represents a decode branch. A PME in MIMD mode enters the SIMD mode upon executing a specific instruction for the transition.

When PMEs transition dynamically between SIMD and MIMD modes, an MIMD mode is entered by execution of a SIMD "write control register" instruction with the appropriate control bit set to a "1". At the completion of the SIMD instruction, the PME enters the MIMD mode, enables interrupts, and begins fetching and executing its MIMD instructions from the main store location specified by its general register R0. Interrupts are masked or unmasked depending on the state of interrupt masks when the MIMD control bit is set. The PME returns to SIMD mode either by being externally reinitialized or by executing a MIMD "write control register" instruction with the appropriate control bit set to zero.

#### Data communication paths and control

Returning to Figure 7 it will be seen that each PME has 3 input ports 400, and 3 output ports 480 intended for on-chip communication plus 1 I/O port 410, 480 for off chip communications. Existing technology, rather than the processor idea, requires that the off-chip port be byte wide half duplex. Input ports are connected such that data may be routed from input to memory, or from input AR register 408 to output register 408 via direct 16 bit data path 425. Memory would be the data sink for messages targeted at the PME or for messages that were moved in 'store and forward' mode. Messages that do not target the particular PME are sent directly to the required output port, providing a 'circuit switched' mode, when blocking has not occurred. The PME S/W is charged with performing the routing and determining the selected transmission mode. This makes dynamically selecting between 'circuited switched' and 'store and forward' modes possible. This is also another unique characteristic of the PME design.

Thus, our preferred node has 8 PMEs and each PME has 4 output ports (Left, Right, Vertical, and External). Three of the input ports and three of the output ports are 16-bit wide full duplex point-to-point connections to the other PMEs on the chip. The fourth ports are combined in the preferred embodiment to provide a half duplex point-to-point connection to an off-chip PME. Due to pin and power constraints that we

EP 0 570 729 A2

have imposed to make use of the less dense CMOS we employ, the actual off-chip interface is a byte-wide path which is used to multiplex two halves of the inter-PME data word. With special "zipper" circuitry which provides a dynamic, temporary logical breaking of intermodal rings to allow data to enter or leave an array, these external PME ports provide the APAP external I/O array function.

For data routed to the PME memory, normal DMA is supported such that the PME instruction stream must become involved in the I/O processing only at the beginning and end of messages. Finally, data that is being circuit switched to an internal output port is forwarded without clocking. This permits single cycle data transfers within a chip and detects when chip crossings will occur such that the fastest but still reliable communication can occur. Fast forwarding utilizes forward data paths and backward control paths, all operating in transparent mode. In essence, a source looks through several stages to see the acknowledgments from the PME performing a DMA or off-chip transfer.

As seen by FIGURES 7 and 8 Data on a PME input port may be destined for the local PME, or for a PME further down the ring. Data destined for a PME further down the ring may be stored in the local PME main memory and then forwarded by the local PME towards the target PME (store and forward), or the local input port may be logically connected to a particular local output port (circuit switched) such that the data passes "transparently" through the local PME on its way to the target PME. Local PME software dynamically controls whether or not the local PME is in "store and forward" mode or in "circuit switched" mode for any of the four inputs and four outputs. In circuit switched mode, the PME concurrently processes all functions except the I/O associated with the circuit switch; in store and forward mode the PME suspends all other processing functions to begin the I/O forwarding process.

While data may be stored externally of the array in a shared memory or DASD (with external controller), it may be stored anywhere in the memories provided by PMEs. Input data destined for the local PME or buffered in the local PME during "store and forward" operations is placed into local PME main memory via a direct memory access (address) mechanism associated with each of the input ports. A program interrupt is available to indicate that a message has been loaded into PME main memory. The local PME program interprets header data to determine if the data destined for the local PME is a control message which can be used to set up a circuit-switched path to another PME, or whether it is a message to be forwarded to another PME. Circuit switched paths are controlled by local PME software. A circuit switched path logically couples a PME input path directly to an output path without passing through any intervening buffer storage. Since the output paths between PMEs on the same chip have no intervening buffer storage, data can enter the chip, pass through a number of PMEs on the chip and be loaded into a target PME's main memory in a single clock cycle. Only when a circuit switch combination leaves the chip, is an intermediate buffer storage required. This reduces the effective diameter of the APAP array by a number of unbuffered circuit switched paths. As a result data can be sent from a PME to a target PME in as few clock cycles as there are intervening chips, regardless of the number of PMEs in the path. This kind of routing can be compared to a switched environment in which at each node cycles are required to carry data on to the next node. Each of our nodes has 8 PMEs!

#### Memory and Interrupt Levels

The PME contains 32K by 16 bit 420 dedicated storage words. This storage is completely general and can contain both data and program. In SIMD operations all of memory could be data as is characteristic of other SIMD massively parallel machines. In MIMD modes, the memory is quite normal; but, unlike most massively parallel MIMD machines the memory is on the same chip with the PME and is thus, immediately available. This then eliminates the need for cache-ing and cache coherency techniques characteristic of other massively parallel MIMD machines. In the case for instance of Inmos's chip, only 4K resides on the chip, and external memory interface bus and pins are required. These are eliminated by us.

Low order storage locations are used to provide a set of general purpose registers for each interrupt level. The particular ISA developed for the PME uses short address fields for these register references. Interrupts are utilized to manage processing, I/O activities and SAV specified functions (i.e., a PME in normal processing will switch to an interrupt level when incoming I/O initiates). If the level is not masked, the switch is made by changing a pointer in H/W such that registers are accessed from a new section of low order memory and by swapping a single PC value. This technique permits fast level switching and permits SAV to avoid the normal register save operations and also to save status within the interrupt level registers.

The PME processor operates on one of eight program interrupt levels. Memory addressing permits a partitioning of the lower 576 words of memory among the eight levels of interrupts. 84 of these 576 words of memory are directly addressable by programs executing on any of the eight levels. The other 512 words

are partitioned into eight 64 word segments. Each 64 word segment is directly accessible only by programs executing on its associated interrupt level. Indirect addressing techniques are employed for allowing all programs to access all 32K words of PME memory.

The interrupt levels are assigned to the input ports, the BCI, and to error handling. There is a "normal" level, but there is no "privileged", nor "supervisor" level. A program interrupt causes a context switch in which the contents of the PC program counter, status/control register, and selected general registers are stored in specified main memory locations and new values for these registers are fetched from other specified main memory locations.

The PME data flow discussed with reference to FIGURES 7 and 8, may be amplified by reference to the additional sections below. In a complex system, the PME data flow uses the combination of the chip as an array node with memory, processor and IO which sends and receives messages with the BCI that we replicate as the basic building block of an MMP built with our APAP. The MMP can handle many word lengths.

#### 15 PME Multiple Length Data Flow Processing

The system we describe can perform the operations handled by current processors with the data flow in the PME which is 16 bits wide. This is done by performing operations on data lengths which are multiples of 16 bits. This is accomplished by doing the operation in 16 bit pieces. One may need to know the result of each piece (i.e. was it zero, was there a carry out of the high bits of the sum).

Adding two numbers of 48 bits can be an example of data flow. In this example two numbers of 48 bits (a(0-47) and b(0-47)) are added by performing the following in the hardware:

a(32-47) + b(32-47) → ans(32-47)    step one

- 1) save the carry out of high bit of sum
- 2) remember if partial result was zero or non-zero

a(16-31) + b(16-31) + save carry → ans(16-31)    step two

- 1) save the carry out of high bit of sum
- 2) remember if partial result was zero or non-zero from this result and from previous step; if both are zero remember zero; if either is non-zero remember non-zero

a(0-15) + b(0-15) + saved carry → ans(0-15)    final step

- 1) if this piece is zero and last piece was zero ans is zero
- 2) if this piece is zero and last piece was non-zero ans is non-zero
- 3) if this piece is non-zero ans is positive or negative based on sign of sum (assuming no overflow)
- 4) if carry into sign of ans is not-equal to carry out of sign of answer, ans has wrong sign and result is an overflow (can not properly represent in the available bits)

The length can be extended by repeating the second step in the middle as many times as required. If the length were 32 the second step would not be performed. If the length were greater than 48, step two would be done multiple times. If the length were just 16 the operation in step one, with conditions 3 and 4 of the final step would be done. Extending the length of the operands to multiple lengths of the data flow is a technique having a consequence that the instruction usually takes longer to execute for a narrower data flow. That is, a 32 bit add on a 32 bit data flow only takes one pass through the adder logic, while the same add on a 16 bit data flow takes two passes through the adder logic.

What we have done that is interesting is that in the current implementation of the machine we have single instructions which can perform adds/subtracts/compares/moves on operands of length 1 to 8 words (length is defined as part of the instruction). Individual instructions available to the programmer perform the same kind of operations as shown above for steps one, two, and final (except to the programmer the operand length is longer i.e. 16 to 128 bits). At the bare bones hardware level, we are working on 16 bits at a time, but the programmer thinks s/he's doing 16 to 128 bits at a time.



By using combinations of these instructions, operands of any length can be manipulated by the programmer i.e. two instructions can be used to add two numbers of up to 256 bits in length.

#### PME Processor

Our PME processor is different from modern microprocessors currently utilized for MPP applications. The processor portion differences include:

1. the processor is a fully programmable hardware computer (see the ISA description for an instruction set overview) with:
  - o it has a complete memory module shown in the upper right corner (see FIGURE 8),
  - o it has hardware registers with controls required to emulate separate register sets for each interrupt level (shown in upper left corner),
  - o its ALU has the required registers and controls to permit effective multi-cycle integer and floating point arithmetic,
  - o it has I/O switching paths needed to support packet or circuit switched data movement between PMEs interconnected by point-to-point links shown in the lower right corner.
2. This is our minimalist approach to processor design permitting eight replications of the PME per chip with the CMOS DRAM technology.
3. This processor portion of the PME provides about the minimum dataflow width required to encode a fast Instruction Set Architecture (ISA) - see Tables - which is required to permit effective MIMD or SIMD operation of our MMP.

#### PME Resident Software

The PME is the smallest element of the APAP capable of executing a stored program. It can execute a program which is resident in some external control element and fed to it by the broadcast and control interface (BCI) in SIMD mode or it can execute a program which is resident in its own main memory (MIMD mode). It can dynamically switch between SIMD mode and MIMD mode representing SIMD/MIMD mode duality functions, and also the system can execute these dualities at the same time (SIMMIMD mode). A particular PME can make this dynamic switch by merely setting or resetting a bit in a control register. Since SIMD PME software is actually resident in the external control element, further discussion of this may be found in our discussion of the Array Director and in related applications.

MIMD software is stored into the PME main memory while the PME is in SIMD mode. This is feasible since many of the PMEs will contain identical programs because they will be processing similar data in an asynchronous manner. Here we would note that these programs are not fixed, but they can be modified by loading the MIMD program from an external source during processing of other operations.

Since the PME instruction set architecture represented in the Tables is that of a microcomputer, there are few restrictions with this architecture on the functions which the PME can execute. Essentially each PME can function like a RISC microprocessor. Typical MIMD PME software routines are listed below:

1. Basic control programs for dispatching and prioritizing the various resident routines.
2. Communication software to pass data and control messages among the PMEs. This software would determine when a particular PME would go into/out of the "circuit switched" mode. It performs a "store and forward" function as appropriate. It also initiates, sends, receives, and terminates messages between its own main memory and that of another PME.
3. Interrupt handling software completes the context switch, and responds to an event which has caused the interrupt. These can include fail-safe routines and rerouting or reassignment of PMEs to an array.
4. Routines which implement the extended Instruction Set Architecture which we describe below. These routines perform macro level instructions such as extended precision fixed point arithmetic, floating point arithmetic, vector arithmetic, and the like. This permits not only complex math to be handled but image processing activities for display of image data in multiple dimensions (2d and 3d images) and multimedia processes.
5. Standard mathematical library functions can be included. These can preferably include LINPAK and VPSS routines. Since each PME may be operating on a different element of a vector or matrix, the various PMEs may all be executing different routines or differing portions of the same matrix at one time.
6. Specialized routines for performing scatter/gather or sorting functions which take advantage of the APAP nodal interconnection structure and permit dynamic multi-dimensional routing are provided. The routines effectively take advantage of some amount of synchronization provided among the various PMEs, while permitting asynchronous operations to continue. For sorts, there are sort routines. The

APAP is well suited to a Batch Sort. Because that sort requires extensive calculations to determine particular element to compare versus very short comparison cycles. Program synchronization is managed by the I/O statements. The program allows multiple data elements per PME and very large parallel sorts in quite a straight forward manner.

5 While each PME has its own resident software, the systems made from these microcomputers can execute higher level language processes designed for scalar and parallel machines. Thus the system can execute application programs which have been written for UNIX machines, or those of other operating systems, in high level languages such as Fortran, C, C++, FortranD, and so on.

10 It may be an interesting footnote that our processor concepts use an approach to processor design which is quite old. Perhaps thirty years of use of a similar ISA design has occurred in IBM's military processors. We have been the first to recognize that this kind of design can be used to advantage to leapfrog the dead ended current modern microprocessor design when combined with our total PME design to move the technology to a new path for use in the next century.

15 Although the processor's design characteristics are quite different from other modern microprocessors, similar gate constrained military and aerospace processors have used the design since the '60s. It provides sufficient instructions and registers for straight forward compiler development, and both general and signal processing applications are effectively running on this design. Our design has minimal gate requirements, and IBM has implemented some similar concepts for years when embedded chip designs were needed general purpose processing. Our adoption now of parts of the older ISA design permits use of many utilities and other software vehicles which will enable adoption of our systems at a rapid rate because of the existing base and the knowledge that many programmers have of the design concepts.

#### PME I/O

25 The PME will interface to the broadcast and control interface (BCI) bus by either reading data from the bus into the ALU via the path labeled BCI in FIGURE 8 or by fetching instructions from the bus directly into the decode logic (not shown). The PME powers up in SIMD mode and in that mode reads, decodes and executes instructions until encountering a branch. A broadcast command SIMD mode causes the transition to MIMD with instructions fetched locally. A broadcast PME instruction 'INTERNAL DLOW' reverts the state.

30 PME I/O can be sending data, passing data or receiving data. When sending data, the PME sets the CTL register to connect XMIT to either L, R, V, or X. HW services then pass a block of data from memory to the target via the ALU multiplexer and the XMIT register. This processing interleaves with normal instruction operation. Depending upon application requirements, the block of data transmitted can contain raw data for a predefined PME and/or commands to establish paths. A PME that receives data will store input to memory and interrupt the active lower level processing. The interpretation task at the interrupt level can use the interrupt event to do task synchronization or initiate a transparent I/O operation (when data is addressed elsewhere.) During the transparent I/O operation, the PME is free to continue execution. Its CTL register makes it a bridge. Data will pass through it without gating, and it will remain in that mode until an instruction or the data stream resets CTL. While a PME is passing data it cannot be a data source, but it can be a data sink for another message.

#### PME Broadcast Section

45 This is a chip-to-common control device interface. It can be used by the device that serves as a controller to command I/O, or test and diagnose the complete chip.

Input is word sequences (either instruction or data) that are available to subsets of PMEs. Associated with each word is a code indicating which PMEs will use the word. The BCI will use the word both to limit a PME's access to the interface and to assure that all required PMEs receive data. This serves to adjust the BCI to the asynchronous PME operations. (Even when in SIMD mode PMEs are asynchronous due to I/O and interrupt processing.) The mechanism permits PMEs to be formed into groups which are controlled by interleaved sets of command/data words received over the BCI.

50 Besides delivering data to the PMEs, the BCI accepts request codes from the PME combines them, and transmits the integrated request. This mechanism can be used in several ways. MIMD processes can be initiated in a group of processors that all end with an output signal. The 'AND' of signals triggers the controller to initiate a new process. Applications, in many cases, will not be able to load all required SAW in PME memory. Encoded request to the controller will be used to acquire a SAW overlay from perhaps the host's storage system.

EP 0 570 729 A2

The controller uses a serial scan loop through many chips to acquire information on individual chips or PMEs. These loops initially interface to the BCI but can in the BCI be bridged to individual PMEs.

#### Broadcast and Control Interface

The BCI broadcast and control interface provided on each chip provides a parallel input interface such that data or instructions can be sent to the node. Incoming data is tagged with a subset identifier and the BCI includes the functions required to assure that all PMEs within the node, operating within the subset, are provided the data or instructions. The parallel interface of the BCI serves both as a port to permit data to be broadcast to all PMEs and as the instruction interface during SIMD operations. Satisfying both requirements plus extending those requirements to supporting subset operations is completely unique to this design approach.

Our BCI parallel input interface permits data or instructions to be sent from a control element that is external to the node. The BCI contains "group assignment" registers (see the grouping concepts in our above application entitled GROUPING OF SIMD PICKETS) which are associated with each of the PMEs. Incoming data words are tagged with a group identifier and the BCI includes the functions required to assure that all PMEs within the node which are assigned to the dedicated group are provided the data or instructions. The parallel interface of the BCI serves as both a port to permit data to be broadcast to the PMEs during SIMD operations, and as the PME instruction/immediate operand interface during SIMD operations.

The BCI also provides two serial interfaces. The high speed serial port will provide each PME with the capability to output a limited amount of status information. That data is intended to:

1. signal our Array Director 610 when the PME, e.g. 500, has data that needs to be read or that the PME has completed some operation. It passes a message to the external control element represented by the Array Director.
2. provide activity status such that external test and monitor elements can illustrate the status of the entire system.

The standard serial port permits the external control element means for selectively accessing a specific PME for monitor and control purposes. Data passed over this interface can direct data from the BCI parallel interface to a particular PME register or can select data from a particular PME register and route it to the high speed serial port. These control points allow the external control element to monitor and control individual PMEs during initial power up and diagnostic phases. It permits Array Director to input control data so as to direct the port to particular PME and node internal registers and access points. These registers provide paths such that PME of a node can output data to the Array Director, and these registers permit the Array Director to input data to the units during initial power up and diagnostic phases. Data input to access point can be used to control test and diagnostic operations, i.e. perform single instruction step, stop on compare, break points, etc.

#### Node Topology

Our modified hypercube topology connection is most useful for massively parallel systems, while other less powerful connections can be used with our basic PMEs. Within our initial embodiment of the VLSI chip are eight PMEs with fully distributed PME internal hardware connections. The internal PME to PME chip configuration is a two rings of four PMEs, with each PME also having one connection to a PME in the other ring. For the case of eight PMEs in a VLSI chip this is a three dimensional binary hypercube, however our approach in general does not use hypercube organizations within the chip. Each PME also provides for the escape of one bus. In the initial embodiment the escaped buses form one ring are called +X, +Y, +W and +Z, while those from the other ring are labeled similarly except - (minus).

The specific chip organization is referred to as the node of the array, and a node can be in a cluster of the array. The nodes are connected using +X and +Y into an array, to create a cluster. The dimensionality of the array is arbitrary, and in general greater than two which is the condition required for developing a binary hypercube. The clusters are then connected using +W, +Z into a array of clusters. Again, the dimensionality of the array is arbitrary. The result is the 4-dimensional hypercube of nodes. The extension to a 5-dimensional hypercube requires the usage of a 10 PME node and uses the additional two buses, say +E1 to connect the 4-dimensional hypercube into a vector of hypercubes. We have then shown the pattern of extension to either odd or even radix hypercubes. This modified topology limits the cluster to cluster wiring while supporting the advantages of the hypercube connection.

Our wireability and topology configuration for massively parallel machines has advantages in keeping the X and Y dimensions within our cluster level of packaging, and in distributing the W and Z bus connections to all the neighboring clusters. After implementing the techniques described, the product will be wireable, and manufacturable while maintaining the inherent characteristics of the topology defined.

The node consists of  $k \cdot n$  PME's plus the Broadcast and Control Interface (BCI) section. Here "n" represents the number of dimensions or rings, which characterize the modified hypercube, while "k" represents the number of rings that characterize the node. Although a node can contain k rings it is a characteristic of the concept that only two of those rings may provide escape buses. "n" and "k" is limited in our preferred embodiment, because of the physical chip package to  $N=4$  and  $k=2$ . This limitation is a physical one, and different chips sets will allow other and increased dimensionality in the array. In addition to being a part of the physical chip package, it is our preferred embodiment to provide a grouping of PME's that interconnect a set of rings in a modified hypercube. Each node will have 8 PME's with their PME architecture and ability to perform processing and data router functions. As such, n is the dimensionality of the modified hypercube (see following section), i.e., a 4d modified hypercube's node element would be 8 PME's while a 5d modified hypercube's node would be 10 PME's. For visualization of nodes which we can employ, refer to FIGURE 6, as well as FIGURES 9 and 10 for visualization of interconnections and see FIGURE 11 for a block diagram of each node. FIGURES 15 and 16 elaborate on possible interconnections for an APAP.

It will be noted that the application entitled "METHOD FOR INTERCONNECTING AND SYSTEM OF INTERCONNECTED PROCESSING ELEMENTS" of co-inventor David B. Rolfe, filed in the United States Patent and Trademark office on May 13, 1991, under USSN 07/698,888, described the modified hypercube criteria which can preferably be used in connection with our APAP MIMP.

That application is incorporated by reference and describes a method of interconnecting processing elements in such a way that the number of connections per element can be balanced against the network diameter (worst case path length). This is done by creating a topology that maintains many of the well known and desirable topological properties of hypercubes while improving its flexibility by enumerating the nodes of the network in number systems whose base can be varied. When using a base 2 number system this method creates the hypercube topology. The invention has fewer interconnections than a hypercube, uniform connections and preserves the properties of a hypercube. These properties include: 1) large number of alternate paths, 2) very high aggregate bandwidth, and 3) well understood and existing methods that can be used to map other common problem topologies with the topology of the network. The result is a generalized non-binary hypercube with less density. It will be understood that with the preference we have given to the modified hypercube approach, in some applications a conventional hypercube can be utilized. In connecting nodes, other approaches to a topology could be used; however, the ones we describe herein are believed to be new and an advance, and we prefer the ones we describe.

The Interconnection methods for the modified hypercube topology for interconnecting a plurality of nodes in a network of PME's:

1. defines a sets of integers  $e_1, e_2, e_3, \dots$  such the product of all elements equals the number of PME's in the network called M, while the product of all elements in the set excepting  $e_1$  and  $e_2$  is the number of nodes called N, and the number of elements in the set called m defines the dimensionality of the network n by the relationship  $n = m-2$ .

2. addresses a PME located by a set of indexes  $a_1, a_2, \dots, a_m$ , where each index is the PME's position in the equivalent level of expansion such that the index  $a_i$  is in the range of zero to  $e_i-1$  for i equal to 1, 2, to  $m$ , by the formula

$$\{ \dots (a(m) \cdot e(m-1) + a(m-2))e(m-1) \dots a(2)e(1) \} + a(1)$$

where the notation  $a(i)$  has the normal meaning of the i-th element in the list of elements called a, or equivalently for e.

3. connects two PME's (with addresses f and g) if and only if either of the following two conditions hold:

a. the integer part of  $f/(e_1 \cdot e_2)$  equals the integer part of  $g/(e_1 \cdot e_2)$  and:

1. the remainder part of  $f/e_1$  differs from the remainder part of  $g/e_1$  by 1 or,

2. the remainder part of  $f/e_2$  differs from the remainder part of  $g/e_2$  by 1 or  $e_2-1$ .

b. the remainder part of  $f/e_1$  differs from the remainder part of  $g/e_1$  for i in the range 3, 4, ... m and the

remainder part of  $f/e_1$  equals the remainder part of  $g/e_2$  which equals i minus three, and the remainder part of  $f/e_2$  differs from the remainder part of  $g/e_2$  by  $e_2$  minus one.

As a result the computing system nodes will form a non-binary hypercube, with the potential for being different radix in each dimension. The node is defined as an array of PME's which supports  $2^n$  ports such

EP 0 570 729 A2

that the ports provided by nodes match the dimensionality requirements of the modified hypercube. If the set of integers  $e_3, e_4, \dots, e_n$ , which define the specific extent of each dimension of a particular modified hypercube are all taken as equal, say  $b$ , and if  $e_1$  and  $e_2$  are taken as 1, then the previous formulas for addressability and connections reduce to:

1.  $N = b^n$
  2. addressing a PME as numbers representing the base  $b$  numbering system
  3. connecting two computing elements ( $f$  and  $g$ ) if and only if the address of  $f$  differs from the address of  $g$  in exactly one base  $b$  digit, using the rule that 0 and  $b-1$  are separated by 1.
  4. the number of connections supported by each PME is  $2^n$
- Which is exactly as described in the base application, with the number of communication buses spanning non-adjacent PMEs chosen as zero.

#### Intra-Node PME Interconnections:

PMEs are configured within the node as a 2 by  $n$  array. Each PME is interconnected with its three neighbors (edges wrap only in the second dimension) using a set of input/output ports, thus, providing full duplex communication capability between PMEs. Each PME's external input and output port is connected to node I/O pins. Input and output ports may be connected to share pins for half-duplex communication or to separate pins for full-duplex capability. The interconnections for a 4d modified hypercube node are shown in FIGURES 9 and 10. (Note that where  $n$  is even the node can be considered to be a 2 by 2 by  $n/2$  array.)

FIGURE 9 illustrates the eight processing elements 500, 510, 520, 530, 540, 550, 560, 570 within the node. The PMEs are connected in a binary hypercube communication network. This binary hypercube displays three intra connections between PMEs (501, 511, 521, 531, 541, 551, 561, 571, 580, 591, 592, 593). Communication between the PME is controlled by in and out registers under control of the processing element. This diagram shows the various paths that can be taken to escape I/O out any of the eight directions,  $+w$  525, 565,  $+x$  515, 555,  $+y$  505, 545,  $+z$  535, 575. The communication can be accomplished without storing the data into memory if desired.

It may be noted that while a network switch chip could be employed to connect various cards each having our chip with other chips of the system, it can and should desirably be eliminated. Our inter PME network that we describe as the "4d torus" is the mechanism used for inter PME-communication. A PME can reach any other PME in the array on this interface. (PMEs in between may be Store/Forward or Circuit Switched)

#### Chip Relationships for Interconnections

We have discussed the chip, and FIGURE 11 shows a block diagram of the PME Processor/Memory chip. The chip consists of the following elements each of which will be described in later paragraphs:

1. 8 PMEs each consisting of a 16 bit programmable processor and 32K words of memory (64K bytes).
2. Broadcast Interface (BCI) to permit a controller to operate all or subsets of the PMEs and to accumulate PME requests,
3. Interconnection Levels
  - a. Each PME supports four 8 bit wide inter-PME communication paths. These connect to 3 neighboring PMEs on the chip and 1 off chip PME.
  - b. Broadcast-to-PME busing, which makes data or instructions available.
  - c. Service Request lines that permit any PME to send a code to the controller. The BCI combines the requests and forwards a summary.
  - d. Serial Service loops permit the controller to read all detail about the functional blocks. This level of interconnection extends from the BCI to all PMEs (FIGURE 11 for ease of presentation omits this detail.)

#### Interconnection and Routing.

The MPP will be implemented by replication of the PME. The degree of replication does not affect the interconnection and routing schemes used. FIGURE 6 provides an overview of the network interconnection scheme. The chip contains 8 PMEs with interconnections to their immediate neighbors. This interconnection pattern results in the three dimensional cube structure shown in FIGURE 10. Each of the processors within the cube has a dedicated bidirectional byte port to the chip's pins; we refer to the set of 8 PMEs as a node.

An  $n$  by  $n$  array of nodes is a cluster. Simple bridging between the  $+$  and  $-x$  ports and the  $+$  and  $-y$  ports provide the cluster node interconnections. Here the our preferred chip or node has 8 PMEs, each of which manages a single external port. This distributes the network control function and eliminates a possible bottleneck for ports. Bridging the outer edges makes the cluster into a logical torus. We have considered clusters with  $n=4$  and  $n=8$  and believe that  $n=8$  is the better solution for commercial applications while  $n=4$  is better for military conduction cooled applications. Our concept does not impose an unchangeable cluster size. On the contrary, we anticipate some applications using variations.

An array of clusters results in the 4 dimensional torus or hypercube structure illustrated in FIGURE 10. Bridging between the  $+$  and  $-w$  ports and  $+$  and  $-z$  ports provides the 4d torus interconnections. This results in each node within a cluster connected to an equivalent node in all adjacent clusters. (This provides 64 ports between two adjacent clusters rather than the 8 ports that would result from larger clusters.) As with the cluster size, the scheme does not imply a particular size array. We have considered  $2 \times 1$  arrays desirable for workstations and MIL applications and  $4 \times 4$ ,  $4 \times 8$  and  $8 \times 8$  arrays for mainframe applications.

Developing an array of 4d toruses is beyond the gate, pin, and connector limitations of our current preferred chip. However, that limitation disappears with our alternative on-chip optical driver/receiver is employed. In this embodiment our network could use an optical path per PME; with 12 rather than 8 PMEs per chip the array of 4d toruses with multi-Tflop (Teraflop) performance and it seems to be economically feasible to make such machines available for the workstation environment. Remember that such alternative machines will use the application programs developed for our current preferred embodiment.

#### 4d cluster Organization

For constructing a 4d modified hypercube 360, as illustrated in FIGURES 8 and 10, nodes supporting 8 external ports 315 are required. Consider the external ports to be labeled as  $+X$ ,  $+Y$ ,  $+Z$ ,  $+W$ ,  $-X$ ,  $-Y$ ,  $-Z$ ,  $-W$ . Then using  $m_1$  nodes, a ring can be constructed by connecting the  $+X$  to  $-X$  ports. Again  $m_2$  such rings can be interconnected into a ring of rings by interconnecting the matching  $+Y$  to  $-Y$  ports. This level of structure will be called a cluster 320. With  $m_1 = m_2 = 8$  it provides for 512 PMEs and such a cluster will be a building block for several size systems (330, 340, 350), as illustrated with  $m=8$  in FIGURE 6.

#### 4d Array Organization

For building large fine-grained systems, sets of  $m_2$  clusters are interconnected in rows using the  $+Z$  to  $-Z$  ports. The  $m_1$  rows are then interconnected using the  $+W$  to  $-W$  ports. For  $m_1 = \dots = m_4 = 8$  this results in system with 32768 or  $8^{+1}$  PMEs. The organization does not require that every dimension be equally populated as shown in FIGURE 6 (large fine-grained parallel processor 370). In the case of the fine-grained small processor, only a cluster might be used with the unused  $Z$  and  $W$  ports being interconnected on the card. This technique saves card connector pins and makes possible the application of this scalable processor to workstations 340, 350 and avionics applications 330, both of which are connector pin limited. Connecting  $+/-$  ports together in the  $Z$  and  $W$  pairs leads to a workstation organization that permits debug, test and large machine software development.

Again, much smaller scale versions of the structure can be developed by generating the structure with a value smaller than  $m=8$ . This will permit construction of single card processors compatible with the requirements for accelerators in the PS/2 or FtSC System 6000 workstation 350.

#### IO Performance

I/O performance includes overhead to setup transfers and actual burst rate data movement. Setup overhead depends upon application function I/O complexity and network contention. For example, an application can program circuit switched traffic with buffering to resolve conflicts or it can have all PMEs transmit left and synchronize. In the first case, I/O is a major task and detailed analysis would be used to size it. We estimate that simple case setup overhead is 20 to 30 clock cycles or .8 to 1.2  $\mu$ -sec.

Burst rate I/O is the maximum rate a PME can transfer data (with either an on or off chip neighbor.) Memory access limits set the data rate at 140 nsec per byte, corresponding to 7.14 Mbytes/s. This performance includes buffer address and count processing plus data read/writes. It uses seven 40ns cycles per 16 bit word transferred.

This burst rate performance corresponds to a cluster having a maximum potential transfer rate of 3.66 Gbytes/s. It also means that a set of eight nodes along a row or column of the cluster will achieve 57 Mbytes/s burst data rate using one set of their 8 available ports. This number is significant because I/O with

EP 0 570 729 A2

the external world will be done by logically 'unzipping' an edge of the wrapped cluster and attaching it to the external system bus.

#### Inter-PME Routing Protocol

The SIMD/MIMD PME comprises interprocessor communication to the external I/O facilities, broadcast control interfaces, and switching features which allow both SIMD/MIMD operation within the same PME. Embedded in the PME is the fully distributed programmable I/O router for processor communication and data transfers between PMEs.

The PMEs have fully distributed interprocessor communication hardware to on-chip PMEs as well as to the external I/O facilities which connect to the interconnected PMEs in the modified hypercube configuration. This hardware is complemented with the flexible programmability of the PME to control the I/O activity via software. The programmable I/O router functions provide for generating data packets and packet addresses. With this information the PME can send the information thru the network of PMEs in a directed method or out multiple paths determined by any fault tolerance requirements.

Distributed fault tolerance algorithms or program algorithms can take advantage of the programmability along with the supported circuit switched modes of the PME. This performance combinational mode enables everything from off-line PMEs or optimal path data structures to be accomplished via the programmable I/O router.

Our study of applications reveals that it is sometimes most efficient to send bare data between PMEs. At other times applications require data and routing information. Further, it is sometimes possible to plan communications so that network conflicts cannot occur; other applications offer the potential for deadlock, unless mechanisms for buffering messages at intermediate nodes are provided. Two examples illustrate the extremes. In the relaxation phase of a PDE solution, each grid point can be allocated to a node. The inner loop process of acquiring data from a neighbor can easily be synchronized over all nodes. Alternatively, image transformations use local data parameters to determine communication target or source identifiers. This results in data moves through multiple PMEs, and each PME becomes involved in the routing task for each packet. Preplanning such traffic is generally not possible.

To enable the network to be efficient for all types of transfer requirements, we partition, between the H/W and S/W, the responsibility for data routing between PMEs. S/W does most of the task sequencing function. We added special features to the hardware (H/W) to do the inner loop transfers and minimize software (S/W) overhead on the outer loops.

I/O programs at dedicated interrupt levels manage the network. For most applications, a PME dedicates four interrupt levels to receiving data from the four neighbors. We open a buffer at each level by loading registers at the level, and executing the IN (it uses buffer address and transfer count but does not await data receipt) and RETURN instruction pair. Hardware then accepts words from the particular input bus and stores them to the buffer. The buffer full condition will then generate the interrupt and restore the program counter to the instruction after the RETURN. This approach to interrupt levels permits I/O programs to be written that do not need to test what caused the interrupt. Programs read data, return, and then continue directly into processing the data they read. Transfer overhead is minimized as most situations require little or no register saving. Where an application uses synchronization on I/O, as in the PDE example, then programs can be used to provide that capability.

Write operations can be started in several ways. For the PDE example, at the point where a result is to be sent to a neighbor, the application level program executes a write call. The call provides buffer location, word count and target address. The write subroutine includes the register loads and OUT instructions needed to initiate the H/W and return to the application. H/W does the actual byte by byte data transfer. More complicated output requirements will use an output service function at the highest interrupt level. Both application and interrupt level tasks access that service via a soft interrupt.

Setting up circuit switched paths builds on these simple read and write operations. We start with all PMEs having open buffers sized to accept packet headers but not the data. A PME needing to send data initiates the transfer by sending an address/data block to a neighboring PME whose address buffer matches the target. In the neighboring PME address information will be stored; due to buffer full an interrupt will occur. The interrupt S/W tests the target address and will either extend the buffer to accept the data or write the target address to an output port and set the CTL register for transparent data movement. (This allows the PME to overlap its application executions with the circuit switched bridging operation.) The CTL register goes to busy state and remains transparent until reset by the presence of a signal at end of data stream or abnormally by PME programming. Any number of variations on these themes can be implemented.

### System I/O and Array Director

FIGURE 12 shows an Array Director in the preferred embodiment, which may perform the functions of the controller of FIGURE 13 which describes the system bus to array connections. FIGURE 13 is composed of two parts, (a) the bus to/from a cluster, and part (b) the communication of information on the bus to/from a PME. Loading or unloading the array is done by connecting the edges of clusters to the system bus. Multiple system buses can be supported with multiple clusters. Each cluster supports 50 to 57 Mbyte/s bandwidth. Loading or unloading the parallel array requires moving data between all or a subset of the PMEs and standard buses (ie MicroChannel, VME-bus, FutureBus, etc). Those buses, part of the host processor or array controller, are assumed to be rigidly specified. The PME Array therefore must be adapted to the buses. The PME Array can be matched to the bandwidth of any bus by interleaving bus data onto  $n$  PMEs, with  $n$  picked to permit PMEs both I/O and processing time. FIGURE 13 shows how we might connect the system buses to the PMEs at two edges of a cluster. Such an approach would permit 114 Mbyte/s to be supported. It also permits data to be loaded at half the peak rate to two edges simultaneously. Although this reduces the bandwidth to 57 Mbyte/s/cluster, it has the advantage of providing orthogonal data movement within the array and ability to pass data between two buses. (We use those advantages to provide fast transpose and matrix multiply operation.)

As shown in part (a) of FIGURE 13, the bus "dots" to all paths on the edges of the cluster; and, the controller generates a gate signal to each path in the required interleave timing. If required to connect to a system bus with greater than 57 Mbyte/s, the data will be interleaved over multiple clusters. For example, in a system requiring 200 Mbyte/s system buses, groups of 2 or 4 clusters will be used. A large MPP has the capacity to attach 16 or 64 such buses to its  $xy$  network paths. By using the  $w$  and  $z$  paths in addition to the  $x$  and  $y$  paths, that number could be doubled.

FIGURE 13 part (b) shows how the data routes to individual PMEs. The FIGURE shows one particular  $w, x, y$  or  $z$  path that can be operated at 7.13 Mbyte/s in burst mode. If the data on the system bus occurred in bursts, and if the PME memory could contain the complete burst, then only one PME would be required. We designed the PME I/O structure to require neither of these conditions. Data can be gated into the PMEx0 at the full rate until buffer full occurs. At that instant, PMEx0 will change to transparent and PMEx1 will begin accepting the data. Within PMEx0 processing of the input data buffer can begin. PMEs that have taken data and processed it are limited because they cannot transmit the results while in the transparent mode. The design resolves this by switching the data stream to the opposite end of the path at intervals. FIGURE 13(b) shows that under S/W control one might dedicate PMEx0 through PMEx3 to accepting data while PMEx12 through PMEx15 unload results and visa-versa. The controller counts words and adds end of block signals to the data stream, causing the switch in direction. One count applies to all paths supported by the controller so controller workload is reasonable.

### SYSTEMS FOR ALTERNATIVE COMPUTERS

FIGURE 18 illustrates a system block diagram for a host attached large system with a single application processor interface (API). This illustration may also be viewed with the understanding that our invention may be employed in stand alone system which use multiple application processor interfaces (not shown). This configuration will support DASD/Graphics on all or many clusters. Workstation accelerators can eliminate the host, application processor interface (API) and cluster synchronizer (CS) illustrated by emulation. The CS not always required. It will depend on the type of processing that is being performed, as well as the physical drive or power provided for a particular application which uses our invention. An application this is doing mostly MMIO processing will not place as high a workload demand on the controller, so here the control bus can see very slow pulse rise times. Conversely, system doing mostly asynchronous A-SIMD operations with many independent groupings may require faster control busing. In this case, a cluster synchronizer will be desirable.

The system block diagram of FIGURE 18 illustrates that a system might consist of host, array controller and PME array. The PME array is a set of clusters supported by a set of cluster controllers (CC). Although a CC is shown for each cluster that relationship is not strictly required. The actual ratio of clusters to CCs is flexible. The CC provides redrive to, and accumulation from the 64 BCIs/clusters. Therefore, physical parameters can be used establish the maximum ratio. Additionally, the CC will provide for controlling multiple independent subsets of the PME array; that service might also become a gating requirement. A study can be made to determine these requirements for any particular application of our invention. Two versions of the CC will be used. A cluster that is to be connected to a system bus requires the CC providing interleave controls (see System I/O and FIGURE 18) and tri-state drivers. A more simple version that omits



EP 0 570 729 A2

the tri-state busing features can also be employed. In the case of large systems, a second stage of redrive and accumulation is added. This level is the cluster synchronizer (CS). The set of CCs plus CS and the Application Processor Interface (API) make up the Array Controller. Only the API is a programmable unit.

Several variations of this system synthesis scheme will be used. These result in different hardware configurations for various applications, but they do not have a major impact on the supporting software.

For a workstation accelerator, the cluster controllers will be attached directly to the workstation system bus; the function of the API will be performed by the workstation. In the case of a RISC-8000, the system bus is a Micro Channel and the CC units can plug directly into the slots within the workstation. This configuration places the I/O devices (DASD, SCSI and display interfaces) on the same bus that loads/unloads the array. As such the parallel array can be used for I/O intensive tasks like real time image generation or processing. For workstations using other bus systems (VME-bus, FutureBus, etc.), a gateway interface will be used. Such modules are readily available in the commercial marketplace. Note that in these minimal scale systems a single CC can be shared between a determined number of clusters, and neither a CS nor an API is needed.

A MIL avionics application might be similar in size to a workstation, but it needs different interfacing. Consider what may become the normal military situation. An existing platform must be enhanced with additional processing capability, but funding prohibits a complete processing system redesign. For this we would attach to the APAP array a smart memory coprocessor. In this case, a special application program interface API that appears to the host as memory will be provided. Data addressed to the host's memory will then be moved to the array via CC(s). Subsequent writes to memory can be detected and interpreted as commands by the API so that the accelerator appears to be a memory mapped coprocessor.

Large systems can be developed as either host attached or as stand alone configurations. For a host attached system, the configuration shown in FIGURE 18 is useful. The host will be responsible for I/O, and the API would serve as a dispatched task manager. However, a large stand alone system is also possible in special situations. For example, a database search system might eliminate the host, attach DASD to the MicroChannels of every cluster and use multiple APIs as bus masters slaved to the PMEs.

#### Zipper Array Interface with External I/O

Our zipper provides a fast I/O connection scheme and is accomplished by placing a switch between two nodes of the array. This switch will allow for the parallel communication into and out of the array. The fast I/O will be implemented along one edge of the array rings and acts like a large zipper into the X, Y, W, Z rings. The name "zipper connection" is given to the fast I/O. Allowing data to be transferred into and out of the network while only adding switch delays to transfer the data between processors is a unique loading technique. The switching scheme does not disrupt the ring topology created by the X, Y, W, Z buses and special support hardware allows the zipper operation to occur while the PE is processing or routing data.

The ability to bring data into and out of a massively parallel system rapidly is an important enhancement to the performance of the overall system. We believe that the way we implement our fast I/O without reducing the number of processors or dimension of the array network is unique in the field of massively parallel environments.

The modified hypercube arrangement can be extended to permit a topology which comprises rings within rings. To support the interface to the external I/O any or all of the rings can be logically broken. The two ends of the broken ring can then be connected to external I/O buses. Breaking the rings is a logical operation so as to permit regular inter-PME communication at certain time intervals while permitting I/O at other time intervals. This process of breaking a level of rings within the modified hypercube effectively "unzips" rings for I/O purposes. The fast I/O "zipper" provides a separate interface into the array. This zipper may exist on 1 to n edges of the modified hypercube and could support either parallel input into multiple dimensions of the array or broadcast to multiple dimensions of the array. Further data transfers into or out of the array could alternate between the two nodes directly attached to the zipper. This I/O approach is unique and it permits developing different zipper sizes to satisfy particular application requirements. For example, in the particular configuration shown in FIGURE 8, called the large fine-grained processor 380, the zipper for the Z and W buses will be dotted onto the MCA bus. This approach optimizes the matrix transposition time, satisfying a particular application requirement for the processor. For a more detailed explanation of the "zipper" structure, reference may be had to the APAP I/O ZIPPER application filed concurrently herewith. The zipper is here illustrated by Figure 14.

Depending on the configuration and the need of the program to roll data and program into and out of the individual processing elements, the size of the zipper can be varied. The actual speed of the I/O zipper is approximately the number of rings attached times the PME bus width, times the PME clock rate all

divided by 2. (The division permits the receiving PME time to move data onward. Since it can send it to any of  $n$  places I/O contention is completely absorbed over the Array.) With existing technology, i.e., 5 MB/sec PME transfer rate, 64 rings on the zipper, and interleaved to two nodes transfers, 320MB/sec Array transfer rates are possible. (See the typical zipper configuration in FIGURE 14). FIGURE 14 illustrates the fast I/O or the so-called "zipper connection" 700, 710 which exists as a separate interface into the array. This zipper may exist on one 700 or two edges 700, 710 of the hypercube network by dotting onto the broadcast bus 720, 730, 740, 750, at multiple nodes in the array 751, 752, 753, 754 and in multiple directions 770, 780, 790, 751, 755, 757.

Today's MCA bus supports 80 to 160 MB per second burst transfer rate and therefore is a good match for a single zipper in simple or non-interleaved mode. The actual transfer rate given channel overhead and efficiency is something less than that. For systems that have even more demanding I/O requirements, multiple zippers and MCA buses can be utilized. These techniques are seen to be important to processors that would support a large external storage associated with nodes or clusters, as might be characteristic of database machines. Such I/O growth capability is completely unique to this machine and has not previously been incorporated in either massively parallel, conventional single processor, or coarse-grained parallel machines.

#### Array Director Architecture

Our massively parallel system is made up of nodal building blocks of multiprocessor nodes, clusters of nodes, and arrays of PMEs already packaged in clusters. For control of these packaged systems we provide a system array director which with the hardware controllers performs the overall Processing Memory Element (PME) Array Controller functions in the massively parallel processing environment. The Director comprises of three functional areas, the Application Interface, the Cluster Synchronizer, and normally a Cluster Controller. The Array Director will have the overall control of the PME array, using the broadcast bus and our zipper connection to steer data and commands to all of the PMEs. The Array Director functions as a software system interacting with the hardware to perform the role as the shell of the operating system. The Array Director in performing this role receives commands from the application interface and issuing the appropriate array instructions and hardware sequences to accomplish the designated task. The Array Director's main function is to continuously feed the instructions to the PMEs and route data in optimal sequences to keep the traffic at a maximum and collisions to a minimum.

The APAP computer system shown in FIGURE 6 is illustrated in more detail in the diagram of FIGURE 12 which illustrates the Array Director which can function as a controller, or array controller, as illustrated in FIGURE 13 and FIGURES 18 and 19. This Array Director 610 illustrated in FIGURE 12 is shown in the preferred embodiment of an APAP in a typical configuration of  $n$  identical array clusters 685, 670, 680, 690, with an array director 610 for the clusters of 512 PMEs, and an application processor interface 630 for the application processor or processors 600. The synchronizer 650 provides the needed sequences to the array or cluster controller 640 and together they make up the "Array Director" 610. The application processor interface 630 will provide the support for the host processor 600 or processors and test/debug workstations. For APAP units attached to one or more hosts, the Array Director serves as the interface between the user and the array of PMEs. For APAPs functioning as stand alone parallel processing machines, the Array Director becomes the host unit and accordingly becomes involved in unit I/O activities.

The Array Director will consist of the following four functional areas: (see the functional block diagram in FIGURE 12)

1. Application Processor Interface (API) 600,
2. Cluster Synchronizer (CS) 650 (8 x 8 array of clusters),
3. Cluster Controller (CC) 640 (8 x 1 array of nodes),
4. Fast I/O (zipper Connection) 620.

#### The Application Processor Interface (API) 600:

When operating in attached modes, one API will be used for each host. That API will monitor the incoming data stream to determine what are instructions to the Array clusters 685, 670, 680, 690 and what are data for the Fast I/O (zipper) 620. When in standalone mode, the API serves as the primary user program host.

To support these various requirements, the APIs contain the only processors within the Array Director, plus the dedicated storage for the API program and commands. Instructions received from the host can call for execution of API subroutines, loading of API memory with additional functions, or for loading of CC and

EP 0 570 728 A2

PME memory with new S/W. As described in the S/W overview section, these various type requests can be restricted to subset of users via the initial programs loaded into the API. Thus, the operating program loaded will determine the type of support provided which can be tailored to match the performance capability of the API. This further permits the APAP to be adjusted to the needs of multiple users requiring managed and well tested services, or to the individual user wishing to obtain peak performance on a particular application.

The API also provides for managing the path to and from the I/O zipper. Data received from the host system in attached modes, or from devices in standalone modes is forwarded to the Array. Prior to initiating this type of operation the PMEs within the Array which will be managing the I/O are initiated. PMEs operating in MIMD mode can utilize the fast interrupt capability and either standard S/W or special functions for this transfer while those operating in SIMD modes would have to be provided detailed control instructions. Data being sent from the I/O zipper requires somewhat the opposite conditioning. PMEs operating in MIMD modes must signal the API via the high speed serial interface and await a response from the API, while PMEs in SIMD modes are already in synchronization with the API and can therefore immediately output data. The ability to system switch between modes provides a unique ability to adjust the program to the application.

#### Cluster Synchronizer (CS) 650

The CS 650 provides the bridge between the API 630 and CC 640. It stores API 630 output in FIFO stacks and monitors the status being returned from the CC 650 (both parallel input acknowledges and high speed serial bus data) to provide the CC, in timely fashion, with the desired routines or operations that need to be started. The CS provides the capability to support different CCs and different PMEs within clusters so as to permit dividing the array into subsets. This is done by partitioning the array and then commanding the involved cluster controllers to selectively forward the desired operation. The primary function of the synchronizer is to keep all clusters operating and organized such that overhead time is minimized or buried under the covers of PME execution time. We have described how the use of the cluster synchronizer in A-SIMD configurations is especially desirable.

#### Cluster Controller (CC) 640

The CC 640 interfaces to the node Broadcast and Control Interface (BCI) 605 for the set of nodes in an array cluster 835. (For a 4d modified hypercube with 8 nodes per ring that means the CC 640 is attached to 64 BCIs 605 in an 8 by 8 array of nodes and is controlling 512 PMEs. Sixty-four such clusters, also in a 8 by 8 array, lead to the full up system with 32768 PMEs.) The CC 640 will send commands and data supplied by the CS 650 to the BCI parallel port and return the acknowledgement data to the CS 650 when operating in MIMD modes. In SIMD mode the interface operates synchronously, and step by step acknowledgments are not required. The CC 640 also manages and monitors the high speed serial port to determine when PMEs within the nodes are requesting services. Such requests are passed upward to the CS 650 while the raw data from the high speed serial interface is made available to the status display interface. The CC 640 provides the CS 650 with an interface to specific nodes within the cluster via the standard speed serial interface.

In SIMD mode the CC will be directed to send instructions or addresses to all the PMEs over the broadcast bus. The CC can dispatch 18 bit instruction to all PMEs every 40 nanoseconds when in SIMD mode. By broadcasting groups of native instructions to the PME, the emulated instruction set is formed.

When in MIMD mode the CC will wait for the endop signal before issuing new instructions to the PMEs. The concept of the MIMD mode is to build strings of micro-routines with native instructions resident in the PME. These strings can be grouped together to form the emulated instructions, and these emulated instruction can be combined to produce service-canned routines or library functions.

When in SIMD/MIMD (SIMMIMD) mode, the CC will issue instruction as if in SIMD mode and check for endop signals from certain PMEs. The PMEs that are in MIMD will not respond to the broadcast instructions and will continue with their designated operation. The unique status indicators will help the CC to manage this operation and determine when and to whom to present the sequential instructions.

#### Operational Software Levels

This application overviews the operational software S/W levels to provide further explanation of the services performed by various hardware H/W components.

Computer systems generally used have an operating system. Operating system kernels which are relatively complete must be provided in most massive MIMD machines, where workstation class CPU chips run kernels such as Mach. The operating system kernel supports message passing or memory coherency. Other massively parallel systems based upon SIMD models have almost no intelligence in the array. There are no "program counters" out in the array, and thus no programs to execute locally. All instructions are broadcast.

In the systems we have provided with our PME as the basis for cluster arrays, there is not need for an operating system at each chip, a node. We provide a library of key functions for computation and/or communication within each PE (PME) that can be invoked at a high level. SIMD-like instructions are broadcast to the array to set each of a selected set of PMEs. These PMEs can then perform in full MIMD mode one or more of these library routines. In addition, basic interrupt handler and communications routines are resident in each PME allowing the PME to handle communication on a dynamic basis. Unlike existing MIMD machines, the APAP structure need not include an entire program in PME memory. Instead all of that code, which is essentially serial, is the cluster controller. Thus such code, 90% by space and 10% by time (typically) can be broadcast in a SIMD fashion to an array of PMEs. Only the truly parallel inner loops are distributed to the PMEs in a dynamic fashion. These are then initiated into MIMD mode just as other "library" routines are. This enables use of program models which are Single Program Multiple data to be used where the same program is loaded in each PME node, with embedded synchronization code, and executed at the local PME. Design parameters affect bandwidth available on different links, and the system paths are programmatically configurable, allowing high bandwidth links on a target network, and allowing dynamic partition of off chip like PME-to-PME links to provide more bandwidth on specific paths as meets the needs of a particular application. The links leaving a chip mate directly with each other, without the need for external logic. There are sufficient links and there is no predesigned constraint as to which other links they can attach to, so that the system can have a diversity of interconnection topologies, with routing performed dynamically and programmatically.

The system allows usage of existing compilers and parsers to create an executable parallel program which could run on a host or workstation based configuration. Sequential source code for a Single Program Multiple Data system would pass through program analysis, for examination of dependency, data and controls, enabling extension of program source to include call graphs, dependency tables, aliases, usage tables and the like.

Thereafter, program transformation would occur whereby a modified version of the program would be created that extends the degree of parallelism by combining sequences or recognizing patterns to generate explicit compiler directives. A next step would be a data allocation and partitioning step, with message generation, which would analyze data usage patterns and allocate so that elements to be combined would share common indexing, addressing pattern, and these would provide embedded program compiler directives and calls to communication services. At this point the program would pass to a level partitioning step. A level partitioning step would separate the program into portions for execution in ARRAY, in ARRAY CONTROLLER (array director or cluster controller), and HOST. Array portions would be interleaved in sections with any required message passing synchronization functions. At this point, level processing could proceed. Host sources would pass to a level compiler (FORTRAN) for assembly compilation. Controller sources would pass to a microprocessor controller compiler, and items that would be needed by a single PME and not available in a library call would pass to a parser (FORTRAN OR C) to an intermediate level language representation which would generate optimized PME code and Array Controller code. PME code would be created at PME machine level, and would include library extensions, which would pass on load into a PME memory. During execution a PME parallel program in the SPMD process of execution could call upon already coded assembly service functions from a runtime library kernel.

Since the APAP can function as either an attached unit that is closely or loosely coupled with its host or as a stand alone processor, some variation in the upper level S/W models exists. However, these variations serve to integrate the various type applications so as to permit a single set of lower level functions to satisfy all three applications. The explanation will address the attached version S/W first and then the modifications required for standalone modes.

In any system, as illustrated by FIGURE 10, where the APAP is intended to attach to a host processor the user's primary program would exist within the host and would delegate to the APAP unit tasks and associated data as needed to provide desired load balancing. The choice of interpreting the dispatched task's program within the host or the Array Director is a user option. Host level interpretation permits the Array Director to work at interleaving users which do not exploit close control of the Array, while APAP interpretation leads to minimal latency in control branching but tends to limit the APAP time to perform multi-user management tasks. This leads to the concept that the APAP and host can be tightly or loosely

EP 0 570 729 A2

coupled.

Two examples illustrate the extremes:

1. When APAP is attached to 3090 class machines with Floating Point Vector Facilities, user data in compressed form could be stored within the APAP. A host program that called for a vector operation upon two vectors with differing sparseness characteristics would then send instructions to the APAP to realign the data into element by element matching pairs, output the result to the Vector Facility, read answer from the Vector Facility and finally reconfigure data into final sparse data form. Segments of the APAP would be interpreting and building sparse matrix bit maps, while other sections would be calculating how to move data between PMEs such that it would be properly aligned for the zipper.
2. With APAP attached to a small inflight military computer, the APAP could be performing the entire workload associated with Sensor Fusion Processing. The host might initiate the process once, send sensor data as it was received to the APAP and then wait for results. The Array Director would then have to schedule and sequence the PME array through perhaps dozens of processing steps required to perform the process.

The APAP will support three levels of user control:

1. Casual User. S/he works with supplied routines and library function. These routines are maintained at the host or API level and can be evoked by the user via subroutine calls within his program.
2. Customizer User. S/he can write special functions which operate within the API and which directly evoke routines supplied with the API or services supplied with the CC or PME.
3. Development User. S/he generates programs for execution in the CC or PME, depending upon API services for program load and status feedback.

Satisfying these three user levels in either closely or loosely coupled systems leads to the partitioning of H/W control tasks.

#### API Software Tasks

The application program interface API contains S/W services that can test the leading words of data received and can determine whether that data should be interpreted by the API, loaded to some storage within the Array Director or PME, or passed to the I/O zipper.

For data that is to be interpreted, the API determines the required operation and invokes the function. The most common type operation would call for the Array to perform some function which would be executed as a result of API writes to the CS (and indirectly to the CC). The actual data written to the CS/CC would in general be constructed by the API operational routine based upon the parameters passed to the API from the host. Data sent to the CS/CC would in turn be forwarded to the PMEs via the node BCL.

Data could be loaded to either API storage, CC storage, or PME memory. Further, data to be loaded to PME memory could be loaded via either the I/O zipper or via the node BCL. For data to be put into the API memory, the incoming bus would be read then written to storage. Data targeted to the CC memory would be similarly read and then be written to the CC memory. Finally, data for the PME memory (in this case normally new or additional MIMD programs) could be sent to all or selected PMEs via the CS/CC/Node BCL or to a subset of PMEs for selective redistribution via the I/O zipper.

When data is to be sent to the I/O zipper, it could be preceded by inline commands that permit the PME MIMD programs to determine its ultimate target; or, it could be preceded by calls to the API service functions to perform either MIMD initiation or SIMD transmission.

In addition to responding to requests for service received via the host interface, the API program will respond to request from the PMEs. Such requests will be generated on the high speed serial port and will be routed through the CC/CS combination. Requests of this sort can result in the API program's directly servicing the PMEs or accessing the PMEs via the standard speed serial port to determine further qualifying data relative to the service request.

#### PME Software

The software plan includes:

- o Generation of PME resident service routines (that is, 'an extended ISA') for complex operations and I/O management.
- o Definition and development of controller executed subroutines that produce and pass control and parameter data to the PMEs via the BCL bus. These subroutines:
  1. cause a set of PMEs to do mathematical operations on distributed objects,

2. provide I/O data management and synchronization services for PME Array and System Bus interactions.
3. provide startup program load, program overlay and program task management for PMEs.
- o Development of data allocation support services for host level programs, and
- o Development of a programming support system including assembler, simulator, and H/W monitor and debug workstation.

Based upon studies of military sensor fusion, optimization, image transformation, US Post Office optical character recognition and FBI fingerprint matching applications, we have concluded that a parallel processor programmed with vector and array commands (like BLAS calls) would be effective. The underlying programming model must match the PME array characteristics feasible with today's technology. Specifically:

- o PMEs can be independent stored program processors.
- o The array can have thousands of PMEs, and be suitable for fine grained parallelism.
- o Inter-PME networks will have very high aggregate bandwidth and a small 'logical diameter'.
- o Built by network connected microprocessor MIMD standards, each PME is memory limited.

Prior programming on MIMD parallel processors has used task dispatching methodology. Such approaches lead to each PME needing access to an portion of a large program. This characteristic, in combination with the non-shared memory characteristic of the H/W, would exhaust PME memory on any significant problem. We therefore target what we believe is a new programming model, called 'asynchronous SIMD' (A-SIMD) type processing. In this connection see USSN 798,798, filed November 27, 1991 of P. Kogge, which is incorporated herein.

A-SIMD programming in our APAP design means that a group of PMEs will be directed by commands broadcast to them as in SIMD models. The broadcast command will initiate execution of a MIMD function within each PME. That execution can involve data dependent branching and addressing within PMEs, and I/O based synchronization with either other PMEs or the BCI.

Normally, PMEs will complete the processing and synchronize by reading the next command from the BCI.

The A-SIMD approach includes both MIMD and SIMD operating modes. Since the approach imposes no actual time limits on the command execution period, a PME operation that synchronizes on data transfers and executes indefinitely can be initiated. Such functions are very effective in data filtering, DSP, and systolic operations. (They can be ended by either BCI interrupts or by commands over the serial control buses.) SIMD operation results from any A-SIMD control stream that does not include MIMD Mode Commands. Such a control stream can include any of the PMEs native instructions. These instructions are routed directly to the instruction decode logic of the PME. Eliminating the PME instruction fetch provides a higher performance mode for tasks that do not involve data dependent branching.

This programming model (supported by H/W features) extends to permitting the array of PMEs to be divided into independent sections. A separate A-SIMD command stream controls each section. Our application studies show that programs of interest divide into separate phases (i.e. input, input buffering, several processing steps, and output formatting, etc.), suitable for pipeline data processing. Fine-grained parallelism results from applying the n PMEs in a section to a program phase. Applying coarse-grained partitioning to applications often results in discovering small repetitive tasks suitable for MIMD or memory bandwidth limited tasks suitable for SIMD processing. We program the MIMD portions using conventional techniques and program the remaining phases as A-SIMD sections, coded with vectorized commands, sequenced by the array controller. This makes the large controller memory the program store. Varying the number of PMEs per section permits balancing the workload. Varying the dispatched task size permits balancing the BCI bus bandwidth to the control requirements.

The programming model also considers allocating data elements to PMEs. The approach is to distribute data elements evenly over PMEs. In early versions of SAW, this will be done by the programmer or by SAW. We recognize that the IBM parallelizing compiler technologies apply to this problem and we expect to investigate their usage. However, the inter-PME bandwidth provided does tend to reduce the importance of this approach. This links data allocation and I/O mechanism performance.

The H/W requires that the PME initiate data transfers out of its memory, and it supports a controlled write into PME memory without PME program involvement. Input control occurs in the receiving PME by providing an input buffer address and a maximum length. When I/O to a PME results in buffer overflow, H/W will interrupt the receiving PME. The low level I/O functions that will be developed for PMEs build on this service. We will support either movement of raw data between adjacent PMEs or movement of addressed data between any PMEs. The last capability depends upon the circuit switched and store and forward mechanisms. The interpret address and forward operation is important for performance. We have

EP 0 570 729 A2

optimized the H/W and S/W to support the operation. Using one word buffers results in an interrupt upon receipt of address header. Comparing target id with local id permits output path selection. Transfer of the subsequent data words occurs in circuit switched mode. A slight variation on this process using larger buffers results in a store and forward mechanism.

Because of the high performance Inter-PME bandwidth, it is not always necessary or desirable to place data elements within the PME Array carefully. Consider shifting a vector data element distributed across PMEs. Our architecture can send data without an address header, thus, providing for very fast I/O. However, we have found, in many applications, that optimizing a data structure for movement in one direction, penalizes data movement in an orthogonal direction. The penalty in such situations approximates the average cost of randomly routing data in the network. This leads to applications where placing data sequentially or randomly (as opposed to arranging data) results in shorter average process times.

Many applications can be synchronized to take advantage of average access time. (For example, POE relaxation processes acquire data from a neighborhood and thus, can average access over at least four I/O operations.) We believe that after considering the factors applicable to vector and array processes, like scatter/gather or row/column arithmetic, many users will find brute force data allocation to be suitable for the application. However, we know of examples that illustrate application characteristics (like required synchronization or biased utilization of shift directions<sup>1</sup>) that tend to force particular data allocation patterns. This characteristic requires that the tools and techniques developed support either manual tuning of the data placement, or simple and non-optimum data allocation. (We will support the non-optimum data allocation strategy with host level macros to provide near transparent port of vectorized host programs to the MPP. The H/W Monitor workstation will permit the user to investigate the resultant performance.)

FIGURE 19 shows the general S/W development and usage environment. The Host Application Processor is optional in that program execution can be controlled from either the Host or the Monitor. Further, the Monitor will effectively replace the Array Controller in some situations. The environment will support program execution on real or simulated MPP hardware. The Monitor is scenario driven so that the developer doing test and debug operations can create procedures to permit effective operation at any level of abstraction.

FIGURE 20 illustrates the levels of H/W supported within the MPP and the user interfaces to these levels.

We see two potential application programming techniques for the MPP. In the least programmer intensive approach, applications would be written in a vectorized high order language. If the user did not feel that the problem warranted tuning data placement then he would use compile time services to allocate data to the PME Array. The application would use vector calls like BLAS that would be passed to the controller for interpretation and execution on the PME Array. Unique calls would be used to move data between host and PME Array. In summary, the user would not need to be aware of how the MPP organized or processed the data. Two optimization techniques will be supported for this type application:

1. Altering the data allocation by constructing the data allocation table will permit programs to force data placements.
2. Generation of additional vector commands for execution by the array controller will permit tuned subfunctions (ie. calling the Gaussian Elimination as a single operation.)

We also see that the processor can be applied to specialized applications as in those referenced in the beginning of this section. In such cases, code tuned to the application would be used. However, even in such applications the degree of tuning will depend upon how important a particular task is to the application. It is in this situation that we see the need for tasks individually suited to SIMD, MIMD or A-SIMD modes.

These programs will use a combination of:

1. Sequences of PME native instructions passed to an emulator function within the array controller. The emulator will broadcast the instruction and its parameters to the PME set. The PMEs in this SIMD mode will pass the instruction to the decode function, simulating a memory fetch operation.
2. Tight inner loops that can be I/O synchronized will use PME native ISA programs. After initiation from a SIMD mode change, they would run continuously in MIMD mode. (The option to return to SIMD mode via a 'RETURN' instruction exists.)
3. More complicated programs, as would be written in a vectorizing command set, would execute subroutines in the array controller that invoked PME native functions. For example a simplified array controller program to do a BLAS 'SAXPY' command on vectors loaded sequentially across PMEs would

<sup>1</sup> Gaussian Elimination with normal pivoting requires shifting rows but not columns. More than 2:1 performance difference would result from arranging the data such that columns were on the fast shift direction. Even with that there is not an advantage to arranging rows in any particular relationship to the buses.

start sequences within the PMEs that:

- a. Enable PMEs with required x elements via comparison of PME id with broadcast 'inor' and 'X\_addr' values,
- b. Compress the x values via a write to consecutive PMEs,
- 5 c. Calculate the address of PMEs with y elements from broadcast data,
- d. Transmit the compressed x data to the y PMEs,
- e. Do a single precision floating point operation in PMEs receiving x values to complete the operation.

Finally, the SAXPY example illustrates one additional aspect of executing vectorized application programs. The major steps execute in the API and could be programmed by either an optimizer or product developer. Normally, the vectorized application would call rather than include this level o code. These steps would be written as C or Fortran code and will use memory mapped read or writes to control the PME array via the BCI bus. Such a program operates the PME array as a series of MIMD steps synchronized by returns to the API program. Minor steps such as the single precision floating point routines would be developed by the Customizer or Product Developer. These operations will be coded using the native PME 15 ISA and will be tuned to the machine characteristics. In general, this would be the domain of the Product Developer since coding, test and optimization at this level require usage of the complete product development tool set.

The APAP can have applications written in sequential Fortran. The path is quite different. FIGURE 21 outlines a Fortran compiler which can be used. In the first step, it uses a portion of the existing parallelizing 20 compiler to develop program dependencies. The source plus these tables become an input to a process that uses a characterization of the APAP MMP and the source to enhance parallelism.

This MMP is a non-shared memory machine and as such allocates data between the PMEs for local and global memory. The very fast data transfer times and the high network bandwidth reduce the time affect of data allocation, but it still is addressed. Our approach treats part of memory as global and uses a 25 SAW service function. It is also possible to use the dependency information to perform the data allocation in a second alternative. The final step in converting the source to multiple sequential programs is performed by the Level Partitioning step. This partitioning step is analogous to the 'Fortran sup 3:el' work being conducted with DARPA funding. The last process in the compilation is generation of the executable code at all individual functional levels. For the PME this will be done by programming the code generator on an 30 existing compiler system. The Host and API code compilers generate the code targeted to those machines.

The PME can execute MIMD software from its own memory. In general, the multiple PMEs would not be executing totally different programs but rather would be executing the same small program in an asynchronous manner. Three basic types of SAW can be considered although the design approach does not 35 limit the APAP to just these approaches:

1. Specialized emulation functions would make the PME Array emulate the set of services provide by standard user libraries like LINPACK or VPSS. In such an emulation package, the PME Array could be using its multiple set of devices to perform one of the operations required in a normal vector call. This type of emulation, when attached to a vector processing unit, could utilize the vector unit for some operations while performing others internally.
- 40 2. The parallelism of the PME Array could be exploited by operating a set of software that provides a new set of mathematical and service functions in the PMEs. This set of primitives would be the codes exploited by a customizing user to construct his application. The prior example of performing sensor fusion on a APAP attached to a military platform would use such an approach. The customizer would write routines to perform Kalman Filters, Track Optimum Assignment and Threat Assessment using the supplied set of function names. This application would be a series of API call statements, and each call would result in initiating the PME set to perform some basic operation like 'matrix multiply' on data stored within the PME Array.
3. In cases where no effective method, considering performance objectives, or application needs exists then custom SAW could be developed and executed within the PME. A specific example is 'Sort'. Many methods to sort data exist and the objective in all cases is to tune the process and the program to the machine architecture. The modified hypercube is well suited to a Batchier Sort; however, that sort requires extensive calculations to determine particular elements to compare versus very short comparison cycles. The computer program in FIGURE 17 shows a simple example of a PME program 1100 to perform the Batchier Sort 1000 with one element per PME. Each line of the program description would be expanded to 3 to 6 PME machine level instructions, and all PMEs would then execute the program in MIMD mode. Program synchronization is managed via the I/O statements. The program extends to more 55 data elements per PME and to very large parallel sorts in a quite straight forward manner.



EP 0 570 729 A2

CC Storage Contents

Data from the CC storage is used by the PME Array in one of two manners. When the PMEs are operating in SIMD, a series of instructions can be fetched by the CC and passed to the node BCI, thus, reducing load on both the API and CS. Alternatively, functions that are not frequently required, such as PME Fault Reconfiguration S/W, PME Diagnostics, and perhaps conversion routines can be stored in the CC memory. Such functions can then be requested by operating PME MIMD programs or moved to the PMEs at the request of API program directives.

70 Packaging of the 8-Way Modified Hypercube

Our packaging techniques take advantage of the eight PMEs packaged in a single chip and arranged in a N-dimensional modified hypercube configuration. This chip level package or node of the array is the smallest building block in the APAP design. These nodes are then packaged in an 8 X 8 array where the +-X and the +-Y makes rings within the array or cluster and the +-W, and +-Z are brought out to the neighboring clusters. A grouping of clusters make up an array. This step significantly cuts down wire count for data and control for the array. The W and Z buses will connect to the adjacent clusters and form W and Z rings to provide total connectivity around the completed array of various size. The massively parallel system will be comprised of these cluster building blocks to form the massive array of PMEs. The APAP will consist of an 8 X 8 array of clusters, each cluster will have its own controller and all the controllers will be synchronized by our Array Director.

Many trade-offs of wireability and topology have been considered, yet with these considerations we prefer the configuration which we illustrate with this connection. The concept disclosed has the advantage of keeping the X and Y dimensions within a cluster level of packaging, and distributing the W and Z bus connections to all the neighboring clusters.

After implementing the techniques described, the product will be wireable, and manufacturable while maintaining the inherent characteristics of the topology defined.

The concept used here is to mix, match, and modify topologies at different packaging levels to obtain the desired results in terms of wire count. For the method to define the actual degree of modification of the hypercube, refer to the Rolfe modified hypercube patent application referenced above. For the purpose of this preferred embodiment, we will describe two packaging levels to simplify our description. It can be expanded.

The first is the chip design or chip package illustrated by FIGURE 3 and FIGURE 11. There are eight of the processing elements with their associated memory and communication logic encompassed into a single chip which is defined as a node. The internal configuration is classified as a binary hypercube or a 2-degree hypercube where every PME is connected to two neighbors. See the PME-PME communication diagram in FIGURE 9, especially 505, 510, 520, 530, 540, 550, 560, 570.

The second step is that the nodes are configured as an 8 X 8 array to make up a cluster. The fully populated machine is built up of an array of 8 X 8 clusters to provide the maximum capacity of 32768 PMEs. These 4096 nodes are connected in an 8 degree modified hypercube network where the communication between nodes is programmable. This ability to program different routing paths adds flexibility to transmit different length messages. In addition to message length differences, there are algorithm optimizations that can be addressed with these programmability features.

The packaging concept is intended to significantly reduce the off page wire count for each of the clusters. This concept takes a cluster which is defined as a 8 X 8 array of nodes 820, each node 825 having 8 processing elements for a total of 512 PMEs, then to limit the X and Y ring within the cluster and, finally, to bring out the W and Z buses to all clusters. The physical picture could be envisioned as a sphere configuration 800, 810 of 64 smaller spheres 830. See FIGURE 15 for a future packaging picture which illustrates the full up packaging technique, limiting the X and Y rings 800 within the cluster and extending out the W and Z buses to all clusters 810. The physical picture could be envisioned as a sphere configuration of 64 smaller spheres 830.

The actual connection of a single node to the adjacent X and Y neighbors 975 exists within the same cluster. The wiring savings occurs when the Z and W buses are extended to the adjacent neighboring clusters as illustrated in FIGURE 16. Also illustrated in FIGURE 16 is the set of the chips or nodes that can be configured as a sparsely connected 4dimensional hypercube or torus 900, 905, 910, 915. Consider each of the 8 external ports to be labeled as +X, +Y, +Z, +W, -X, -Y, -Z, -W 960, 975. Then, using m chips, a ring can be constructed by connecting the +X to -X ports. Again m such rings can be interconnected into a ring of rings by interconnecting the matching +Y to -Y ports. This level of structure will be called a cluster.

It provides for 512 PME's and will be the building block for several size systems. Two such connections (950, 975) are shown in the diagram illustrated in FIGURE 16.

#### Applications for Deskside MPP.

5

The desktide MPP in a workstation can be effectively applied in several application areas including:

1. Small production tasks that depend upon compute intensive processes. The US Postal Service requires a processor that can accept a fax image of a machine printed envelope and then find and read the zip code. The process is needed at all regional sort facilities and is an example of a very repetitive but still compute intensive process. We have implemented APL language versions of a sample of the required programs. These models emulate the vector and array processes that will be used to do the work on the MPP. Based upon this test, we know that the task is an excellent match to the processing architecture.
2. Tasks in which an analyst, as a result of prior output, or expected needs requests sequences of data transformations. In an example taken from the Defense Mapping Agency, satellite images are to be transformed and smoothed pixel by pixel into some other coordinate system. In such a situation, the transformation parameters for the image vary across localities as a result of ground elevation and slope. The analyst must therefore add fixed control points and reprocess transformations. A similar need occurs in the utilization of scientific simulation results when users require almost real time rotation or perspective changes.
3. Program development for production versions of the MPP will use workstation size MPPs. Consider a tuning process that requires analysis of processor versus network performance. Such a task is machine and analyst interactive. It can require hours when the machine is idle and the analyst is working. When performed on a supercomputer it would be very costly. However, providing an affordable workstation MPP with the same (but scaled) characteristics as the supercomputer MPP eliminates costs and eases the test and debug process by eliminating the programmer inefficiencies related to accessing remote processors.

FIGURE 22 is a drawing of the workstation accelerator. It uses the same size enclosure as the RISC/6000 model 530. Two swing out gates, each containing a full cluster are shown. The combined two clusters provide 5 GOPS of fixed point performance and 530 MflopS of processing power and about 100 Mbytes of I/O bandwidth to the array. The unit would be suitable for any of the prior applications. With quantity production and including a host RISC/6000, it would be price comparable with high performance workstations, not at the price of comparable machines employing old technology.

#### Description of the AWACS Sensor Fusion

The military environment provides a series of examples showing the need for a hardened compute intensive processor.

Communication in the targeted noisy environments implies the need for digitally encoded communications, as is used in ICNIA systems. The process of encoding the data for transmission and recovering information after receipt is a compute intensive process. The task can be done with specialized signal processing modules, but for situations where communication encoding represents bursts of activity, specialized modules are mostly idle. Using the MPP permits several such tasks to be allocated to a single module and saves weight, power, volume and cost.

Sensor data fusion presents a particularly clear example of enhancing an existing platform with the compute power gained from the addition of MPP. On the Air Force E3 AWACS there are more than four sensors on the platform, but there is currently no way to generate tracks resulting from the integration of all available data. Further, the existing generated tracks have quite poor quality due to sampling characteristics. Therefore, there is motivation to use fusion to provide an effective higher sample rate.

We have studied this sensor fusion problem in detail and can propose a verifiable and effective solution, but that solution would overwhelm the compute power available in an AWACS data processor. FIGURE 23 shows the traditional track fusion process. The process is faulty because each of the individual processes tends to make some errors and the final merge tends to collect them instead of eliminating them. The process is also characterized by high time latency in that merging does not complete until the slowest sensor completes. FIGURE 24 presents an improvement and the resulting compute intensive problem with the approach. Although we cannot solve a NP-Hard problem, we have developed a good method to approximate the solution. While the details of that application are being described by the inventors elsewhere, as it can be employed on a variety of machines like an Intel Touchstone with 512 i860 (80960).

EP 0 570 729 A2

processors, and IBM's Scientific Visualization System, it can be used as an application suitable for the MHP using the APAP design described here with say 128,000 PMEs, substantially outperforming these other systems. Application experiments show the approximation quality is below the level of sensor noise and as such the answer is applicable to applications like AWACS. FIGURE 25 shows the processing loop involved in the proposed Lagrangean Reduction n-dimensional Assignment algorithm. The problem uses very controlled repetitions of the well known 2-dimensional assignment problem, the same algorithm that classical sensor fusion processing uses.

Suppose for example that the n-dimensional algorithm was to be applied to the seven sets of observations illustrated in FIGURE 24 and further, suppose that each pass through a reduction process required four iterations through a 2d Assignment process. Then the new 8d Assignment Problem would require 4000 iterations of the 2d Assignment Problem. The AWACS' workload is now about 90% of machine capacity. Fusion perhaps requires 10% of the total effort, but even that small effort when scaled up 4000 times results in total utilization being 370 times the capacity of an AWACS. Not only does this workload overwhelm the existing processor, but it would be marginal in any new MIL environment suited, coarse-grained, parallel processing system currently existing or anticipated in the next few years. If the algorithm required an average of 5 rather than 4 iterations per step, then it would overwhelm even the hypothesized systems. Conversely, the MPP solution can provide the compute power and can do so even at the 6 iteration level.

#### 20 Mechanical Packaging

As illustrated in FIGURE 3, and other FIGURES, our preferred chip is configured in a quadflatpack form. As such it can be brickwalled into various 2 D and 3 D configurations in a package. One chip of eight or more processor memory elements is a first level package module, the same as a single DRAM memory chip is to a foundry which packages the chip. However, it is in a quadflatpack form, allowing connections to one another in four directions. Each connection is point to point. (One chip in its first level package is a module to the foundry.) We are able to construct PE arrays of sufficient magnitude to hit our performance goals due to this feature. The reality is that you can connect these chips across 3, 4 or even five feet, point-to-point, i.e. multi-processor node to node, and still have proper control without the need of fiber optics.

This has an advantage for the drive/receive circuits that are required on the modules. One can achieve high performance and keep the power dissipation down because we do not have bus systems that daisy chain from module to module. We broadcast from node to node, but this need not be a high performance path. Most data operations can be conducted in a node, so data path requirements are reduced. Our broadcast path is essentially primarily used as a controller routing tool. The data stream attaches to and runs in, the ZWXY communication path system.

Our power dissipation is 2.2 watts per node module for our commercial workstation. This allows us to use air cooled packaging. The power system requirements for our system are also reasonable because of this fact. Our power system illustrated multiplies the number of modules supported by about 2.5 watts per module, and such a five volt power supply is very cost effective. Those concerned with the amount of electricity consumed would be astonished that 32 microcomputers could operate with less than the wattage consumed by a reading light.

Our thermal design is enhanced because of the packaging. We avoid hot spots due to high dissipating parts mixed with low dissipating ones. This reflects directly on the cost of the assemblies.

The cost of our system is very attractive compared to the approaches that put a superscalar processor on a card. Our performance level per assembly per watt per connector per part type per dollar is excellent.

Furthermore, we do not need the same number of packaging levels that the other technology does. We do not need module/card/backplane and cable. We can skip the card level if we want to. As illustrated in our workstation modules, we have skipped the card level with our brickwalled approach.

Furthermore, as we illustrated in our layout, each node housing which is brickwalled in the workstation modules, can as illustrated in FIGURE 3 comprise multiple replicated dies, even within the same chip housing. While normally we would place one die within an air cooled package, it is possible to place 8 die on a substrate using a multiple chip module approach. Thus, the envisioned watch with 32 or more processors, is possible, as well as many other applications. The packaging and power and flexibility make applications which are endless. A house could have its controllable instruments all watched, and coordinated with a very small part. Those many chips that are spread around an automobile for engine watching, brake adjustment, and so on could all have a monitor within a housing. In addition, on the same substrate with hybrid technology, one could mount a 386 microprocessor chip with full programmable capability and memory (all in one chip) and use it as the array controller for the substrate package.

We have shown many configurations of systems, from control systems, FIGURE 3, to larger and larger systems. The ability to package a chip with multiple processor memory element of eight or more on a chip in a dip, with pinouts fitting in a standard DRAM memory module, such as in a SIM module make possible countless additional applications ranging from controls to wall size video displays which can have a repetition rate, not a the 15 or so frames that press the existing technology today, but at 30 frames, with a processor assigned to monitor a pixel, or a node only a few pixels. Our brickwall quadflatpack makes it easy to replicate the same part time over and over again. Furthermore, the replicated processor is really memory with processor interchange. Part of the memory can be assigned to a specific monitoring task, and another part (with a size programmatically defined) can be a massive global memory, addressed point-to-point, with broadcast to all capability.

Our basic workstation, our supercomputer, our controller, our AWACS, all are examples of packages that can employ our new technology. An array of memory, with inbuilt CPU chips and I/O, functions as a PME of massively parallel applications, and even more limited applications. The flexibility of packaging and programming makes imaginations expand and our technology allows one part to be assigned to many ideas and images.

#### Military Avionics Applications

The cost advantage of constructing a MIL MPP is particularly well illustrated by the AWACS. It is a 20 year old enclosure that has grown empty space as new technology memory modules have replaced the original core memories. FIGURE 26 shows a MIL qualifiable two cluster system that would fit directly into the rack's empty space and would use the existing memory bus system for interconnection.

Although the AWACS example is very advantageous due to the existence of empty space, in other systems it is possible to create space. Replacing existing memory with a small MPP or gateway to an an isolated MPP is normally quite viable. In such cases, a quarter cluster and a adapter module would result in a 8 Megabyte memory plus 640 MIPs and use perhaps two slots.

#### Supercomputer Application

A 64 cluster MPP is a 13.6 Gflop supercomputer. It can be configured in a system described in FIGURE 27.

The system we describe allows node chips to be brick walled on cluster cards as illustrated in FIGURE 27 to build up systems with some significant cost and size advantages. There is no need to include extra chips such as a network switch in such a system because it would increase costs.

Our interconnection system with "brickwalled" chips allows systems to be built like massive DRAM memory is packaged and will have a defined bus adapter conforming to the rigid bus specifications, for instance a microchannel bus adaptor. Each system will have a smaller power supply system and cooling design than other systems based upon many modern microprocessors.

Unlike most supercomputers our current preferred APAP with floating point emulation is much faster in integer arithmetic (164 GIPS) than it is when doing floating point arithmetic. As such, the processor would be most effective when used in applications that are very character or integer intensive. We have considered three program challenges which in addition to the other applications discussed herein are needful of solution. The applications which may be more important than some of the "grand challenges" to day to day life include:

1. 3000 Vector Processors contain a very high performance floating point arithmetic unit. That unit, as do most vectorized floating point units, requires pipeline operations on dense vectors. Applications that make extensive use of non-regular sparse matrices (i.e. matrices described by bit maps rather than diagonals) waste the performance capability of the floating point unit. The MPP solves this problem by providing the storage for the data and using its compute power and network bandwidth, not to do the calculation but rather to construct dense vectors, and to decompress dense results. The Vector Processing Unit is kept busy by a continual flow of operations on dense vectors being supplied to it by the MPP. By sizing the MPP so that it can effectively compress and decompress at the same rate the Vector Facility processes, one could keep both units fully busy.
2. Another host attached system we considered is a solution to the FBI fingerprint matching problem. Here, a machine with more than 64 clusters was considered. The problem was to match about 8000 fingerprints per hour against the entire database of fingerprint history. Using massive DASD and the full bandwidth of the MPP to host attachment, one can roll the complete data base across the incoming prints in about 20 minutes. Operating about 75% of the MPP in a SIMD mode coarse matching

operation, balances processing to required throughput rate. We estimate that 15% of the machines in A-SIMD processing mode would then complete the matching by doing the detailed verification of unknown print versus file print for cases passing the coarse filter operation. The remaining portions of the machine were in MIMD mode and allocated to reserve capacity, work queue management and output formatting.

3. Application of the MPP to database operations has been considered. Although the work is very preliminary, it does seem to be a good match. Two aspects of the MPP support this premise:

a. The connection between a cluster Controller and the Application Processor Interface is a Micro-Channel. As such, it could be populated with DASD dedicated to the cluster and accessed directly from the cluster. A 64 cluster system with six 640 Mbyte hard drives attached per cluster would provide 246 Gbyte storage. Further, that entire database could be searched sequentially in 10 to 20 seconds.

b. Databases are generally not searched sequentially. Instead they use many levels of pointers. Indexing of databases can be done within the cluster. Each bank of DASD would be supported by 2.5 GIPS of processing power and 32 Mbyte of storage. That is sufficient for both searching and storing the indices. Since indices are now frequently stored within the DASD, significant performance gains would occur. Using such an approach and dispersing DASD on SCSI interfaces attached to the cluster MicroChannel permits effectively unlimited size data bases.

FIGURE 27 is an illustration of the APAP when used to build the system into a supercomputer scaled MPP. The approach reverts to replicating units, but here it is enclosures containing 16 clusters that are replicated. The particular advantage of this replication approach is that the system can be scaled to suit the user's needs.

#### System Architecture

An advantage of the system architecture which is employed in the current preferred embodiment is the ISA system which will be understood by many who will form a pool for programming the APAP. The PME ISA consists of the following Data and Instruction Formats, illustrated in the Tables.

#### Data Formats

The basic (operand) size is the 16 bit word. In PME storage, operands are located on integral word boundaries. In addition to the word operand size, other operand sizes are available in multiples of 16 bits to support additional functions.

Within any of the operand lengths, the bit positions of the operand are consecutively numbered from left to right starting with the number 0. Reference to high-order or most-significant bits always refer to the left-most bit positions. Reference to the low-order or least-significant bits always refer to the right-most bit positions.

#### Instruction Formats

The length of an instruction format may either be 16 bits or 32 bits. In PME storage, instructions must be located on a 16 bit boundary.

The following general instruction formats are used. Normally, the first four bits of an instruction define the operation code and are referred to as the OP bits. In some cases, additional bits are required to extend the definition of the operation or to define unique conditions which apply to the instruction. These bits are referred to as OPX bits.

## EP 0 570 729 A2

| Format Code | Operation            |
|-------------|----------------------|
| RR          | Register to Register |
| DA          | Direct Address       |
| RS          | Register Storage     |
| RI          | Register Immediate   |
| SS          | Storage to Storage   |
| SPC         | Special              |

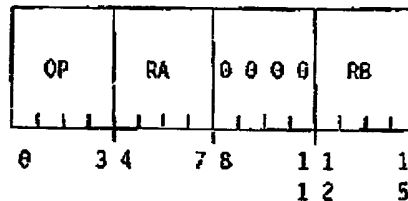
All formats have one field in common. This field and its interpretation is:

**Bits 0-3 Operation Code** - This field, sometimes in conjunction with an operation code extension field, defines the operation to be performed.

Detailed figures of the individual formats along with interpretations of their fields are provided in the following subsections. For some instructions, two formats may be combined to form variations on the instruction. These primarily involve the addressing mode for the instruction. As an example a storage to storage instruction may have a form which involves direct addressing or register addressing.

### RR Format

The Register-Register (RR) format provides two general register addresses and is 16 bits in length as shown.



In addition to an Operation Code field, the RR format contains:

**Bits 4-7 Register Address 1** - The RA field is used to specify which of the 16 general registers is to be used as an operand and/or destination.

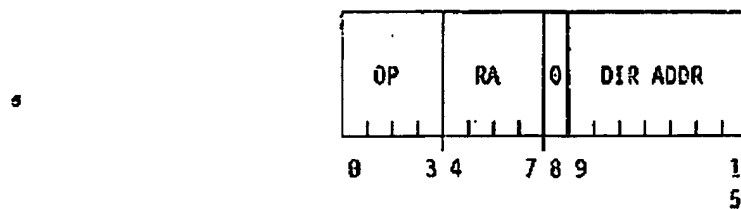
**Bits 8-11 Zeros** - Bit 8 being a zero defines the format to be a RR or DA format and bits 9-11 equal to zero define the operation to be a register to register operation (a special case of the Direct Address format).

**Bits 12-15 Register Address 2** - The RB field is used to specify which of the 16 general registers is to be used as an operand.

### DA Format

The Direct Address (DA) format provides one general register address and one direct storage address as shown.

EP 0 570 729 A2



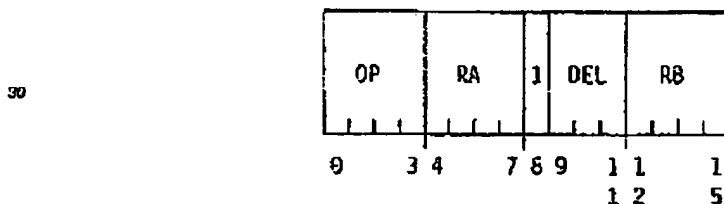
In addition to an Operation Code field, the DA format contains:

- Bits 4-7 **Register Address 1** - The RA field is used to specify which of the 16 general registers is to be used as an operand and/or destination.
  - 15 **Bit 8** **Zero** - This bit being zero defines the operation to be a direct address operation or a register to register operation.
  - Bits 9-15** **Direct Storage Address** - The Direct Storage Address field is used as an address into the level unique storage block or the common storage block. Bits 9-11 of the direct address field must be non-zero to define the direct address form.
- 20

## RS Format

The Register Storage (RS) format provides one general register addresses and an indirect storage address.

25



In addition to an Operation Code field, the RS format contains:

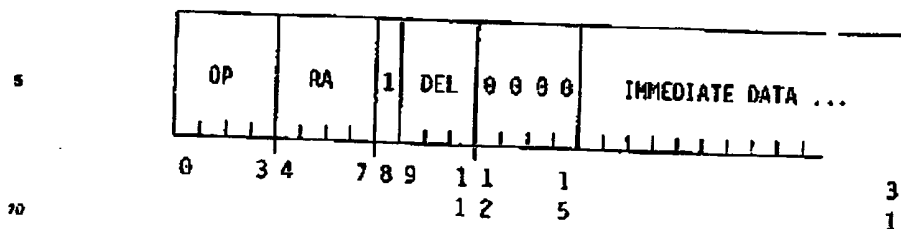
- Bits 4-7 **Register Address 1** - The RA field is used to specify which of the 16 general registers is to be used as an operand and/or destination.
  - 40 **Bit 8** **One** - This bit being one defines the operation to be a register storage operation.
  - Bits 9-11** **Register Data** - These bits are considered a signed value which is used to modify the contents of register specified by the RB field.
  - Bits 12-15** **Register Address 2** - The RB field is used to specify which of the 16 general registers is to be used as an storage address for an operand.
- 45

## RI Format

The Register-Immediate (RI) format provides one general register address and 16 bits of immediate data. The RI format is 32 bits of length as shown:

50

55



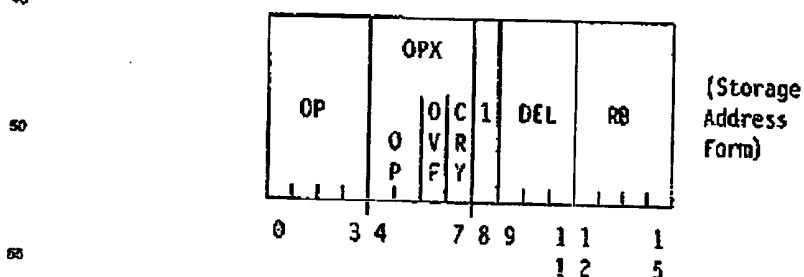
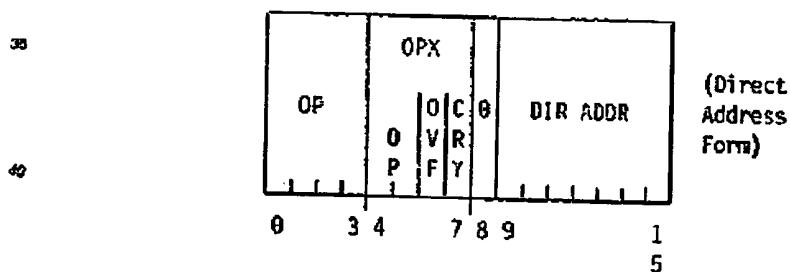
In addition to an Operation Code field, the RI format contains:

- 15
- Bits 4-7** Register Address 1 - The RA field is used to specify which of the 16 general registers is to be used as an operand and/or destination.
- Bit 8** One - This bit being one defines the operation to be a register storage operation.
- Bits 9-11** Register Data - These bits are considered a signed value which is used to modify the contents of the program counter. Normally, this field would have a value of one for the register immediate format.
- 20 **Bits 12-15** Zeros - The field being zero is used to specify that the updated program counter, which points to the immediate data field, is to be used as a storage address for an operand.
- Bits 16-31** Immediate Data - This field serves as a 16 bit immediate data operand for Register Immediate instructions.

25 SS Format

The Storage to Storage (SS) format provides two storage addresses, one explicit; and the second implicit. The implied storage address is contained in General Register 1. Register 1 is modified during execution of the instruction. There are two forms of a SS instruction, a direct address form and a storage

30 address form.



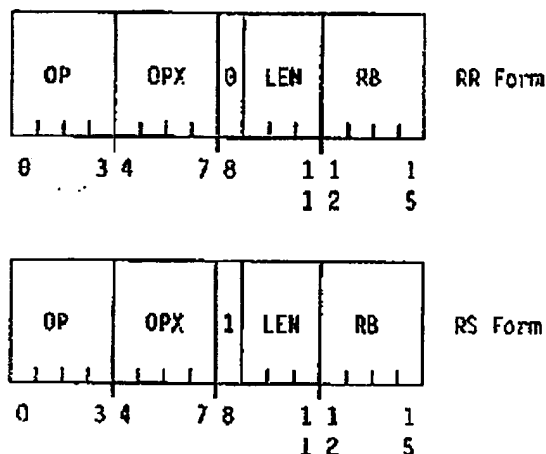


In addition to an Operation Code field, the S8 format contains:

- Bits 4-7 **Operation Extension Code** - The OPX field, together with the Operation Code, defines the operation to be performed. Bits 4-5 define the operation type such as ADD or SUBTRACT. Bits 6-7 control the carry, overflow, and how the condition code will be set. Bit 6 = 0 ignores overflow, bit 6 = 1 allows overflow. Bit 7 = 0 ignore the carry stat during the operation; bit 7 = 1 includes the carry stat during the operation.
- Bit 8 **Zero** - Defines the form to be a direct address form. One - Defines the form to be a storage address form.
- Bits 9-15 **Direct Address (Direct Address Form)** - The Direct Storage Address field is used as an address into the level unique storage block or the common storage block. Bits 9-11 of the direct address field must be non-zero to define the direct address form.
- Bits 9-11 **Register Delta (Storage Address Form)** - These bits are considered a signed value which is used to modify the contents of register specified by the RB field.
- Bits 12-15 **Register Address 2 (Storage Address Form)** - The RB field is used to specify which of the 16 general registers is to be used as a storage address for an operand.

#### SPC Format 1

The Special (SPC1) format provides one general register storage operand address.



In addition to an Operation Code field, the SPC1 format contains:

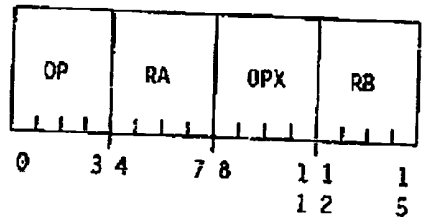
- Bits 4-7 **OP Extension** - The OPX field is used to extend the operation code.
- Bit 8 **Zero or One** - This bit being zero defines the operation to be a register operation. This bit being one defines the operation to be a register storage operation.
- Bits 9-11 **Operation Length** - These bits are considered an unsigned value which is used to specify the length of the operand in 16 bit words. A value of zero corresponds to a length of one, and a value of B'111' corresponds to a length of eight.
- Bits 12-15 **Register Address 2** - The RB field is used to specify which of the 16 general registers is to be used as a storage address for the operand.

#### SPC Format 2

The Special (SPC2) format provides one general register storage operand address.

EP 0 570 729 A2

5



10

In addition to an Operation Code field, the SPC2 format contains:

- Bits 4-7 Register Address 1 - The RA field is used to specify which of the 16 general registers is to be used as an operand and/or destination.
- Bits 8-11 OP Extension - The OPX field is used to extend the operation code.
- Bits 12-15 Register Address 2 - The RB field is used to specify which of the 16 general registers is to be used as a storage address for the operand.

THE INSTRUCTION LIST OF THE ISA INCLUDES THE FOLLOWING:

20

Table 1 (Page 1 of 3). Fixed-Point Arithmetic Instructions

NAME

MNE-

TYPE

25

ADD DIRECT

MONIC

ada

DA

30

35

40

45

50

55

Table 1 (Page 2 of 3). Fixed-Point Arithmetic Instructions

| <u>NAME</u>                               | <u>MNE-</u>  | <u>TYPME</u> |
|---|--------------|--------------|
|   | <u>MONIC</u> |              |
| 5     ADD FROM STORAGE                    | a            | RS           |
| (WITH DELTA)                              | awd          | RS           |
| ADD IMMEDIATE                             | ai           | RI           |
| 10       (WITH DELTA)                     | aiwd         | RI           |
| ADD REGISTER                              | ar           | RR           |
| COMPARE DIRECT ADDRESS                    | cda          | DA           |
| 15     COMPARE IMMEDIATE                  | ci           | RI           |
| (WITH DELTA)                              | ciwd         | RI           |
| COMPARE FROM STORAGE                      | c            | RS           |
| 20       (WITH DELTA)                     | cwd          | RS           |
| COMPARE REGISTER                          | cr           | RR           |
| COPY                                      | cpy          | RS           |
| (WITH DELTA)                              | cpywd        | RS           |
| 25     COPY WITH BOTH IMMEDIATE           | cpybi        | RI           |
| (WITH DELTA)                              | cpybiwd      | RI           |
| COPY IMMEDIATE                            | cpyi         | RI           |
| 30       (WITH DELTA)                     | cpyiwd       | RI           |
| COPY DIRECT                               | cpyda        | DA           |
| COPY DIRECT IMMEDIATE                     | cpydai       | DA           |
| 35     INCREMENT                          | inc          | RS           |
| (WITH DELTA)                              | incwd        | RS           |
| LOAD DIRECT                               | lda          | DA           |
| 40     LOAD FROM STORAGE                  | l            | RS           |
| (WITH DELTA)                              | lwd          | RS           |
| LOAD IMMEDIATE                            | li           | RI           |
| 45       (WITH DELTA)                     | liwd         | RI           |
| LOAD REGISTER                             | lr           | RR           |
| MULTIPLY SIGNED                           | mpy          | SPC          |
| MULTIPLY SIGNED EXTENDED                  | mpyx         | SPC          |
| 50     MULTIPLY SIGNED EXTENDED IMMEDIATE | mpyxi        | SPC          |

55

Table 1 (Page 3 of 3). Fixed-Point Arithmetic Instructions

| NAME                                 | MNE-         | TYPE |
|--------------------------------------|--------------|------|
|                                      | <u>MONIC</u> |      |
| MULTIPLY SIGNED IMMEDIATE            | mpyi         | SPC  |
| MULTIPLY UNSIGNED                    | mpyu         | SPC  |
| MULTIPLY UNSIGNED EXTENDED           | mpyux        | SPC  |
| MULTIPLY UNSIGNED EXTENDED IMMEDIATE | mpyuxi       | SPC  |
| MULTIPLY UNSIGNED IMMEDIATE          | mpyui        | SPC  |
| STORE DIRECT                         | stda         | DA   |
| STORE                                | st           | RS   |
| (WITH DELTA)                         | stwd         | RS   |
| STORE IMMEDIATE                      | sti          | RI   |
| (WITH DELTA)                         | stwd         | RI   |
| SUBTRACT DIRECT                      | sda          | DA   |
| SUBTRACT FROM STORAGE                | s            | RS   |
| (WITH DELTA)                         | swd          | RS   |
| SUBTRACT IMMEDIATE                   | si           | RI   |
| (WITH DELTA)                         | stwd         | RI   |
| SUBTRACT REGISTER                    | sr           | RR   |
| SWAP AND EXCLUSIVE OR WITH STORAGE   | swapx        | RR   |

Table 2 (Page 1 of 3). Storage to Storage Instructions

| NAME                                | MNE-         | TYPE |
|-------------------------------------|--------------|------|
|                                     | <u>MONIC</u> |      |
| ADD STORAGE TO STORAGE              | sa           | SS   |
| (WITH DELTA)                        | sawd         | SS   |
| ADD STORAGE TO STORAGE DIRECT       | sada         | SS   |
| ADD STORAGE TO STORAGE FINAL        | saf          | SS   |
| (WITH DELTA)                        | safwd        | SS   |
| ADD STORAGE TO STORAGE FINAL DIRECT | saida        | SS   |
| ADD STORAGE TO STORAGE INTERMEDIATE | sai          | SS   |
| (WITH DELTA)                        | saiwd        | SS   |

Table 2 (Page 2 of 3). Storage to Storage Instructions

| <u>NAME</u>                                       | <u>MNE-<br/>MONIC</u> | <u>TYPME</u> |
|---|-----------------------|--------------|
| 5 ADD STORAGE TO STORAGE INTERMEDIATE<br>DIRECT   | saida                 | SS           |
| ADD STORAGE TO STORAGE LOGICAL                    | sai                   | SS           |
| 10 (WITH DELTA)                                   | saiwd                 | SS           |
| ADD STORAGE TO STORAGE LOGICAL DIRECT             | saida                 | SS           |
| COMPARE STORAGE TO STORAGE                        | sc                    | SS           |
| 15 (WITH DELTA)                                   | scwd                  | SS           |
| COMPARE STORAGE TO STORAGE DIRECT                 | scda                  | SS           |
| COMPARE STORAGE TO STORAGE FINAL                  | scf                   | SS           |
| 20 (WITH DELTA)                                   | schwd                 | SS           |
| COMPARE STORAGE TO STORAGE FINAL DIRECT           | scfda                 | SS           |
| COMPARE STORAGE TO STORAGE INTERMEDIATE           | sci                   | SS           |
| 25 (WITH DELTA)                                   | sciwd                 | SS           |
| COMPARE STORAGE TO STORAGE INTERMEDIATE<br>DIRECT | scida                 | SS           |
| COMPARE STORAGE TO STORAGE LOGICAL                | scl                   | SS           |
| 30 (WITH DELTA)                                   | sclwd                 | SS           |
| COMPARE STORAGE TO STORAGE LOGICAL<br>DIRECT      | sclda                 | SS           |
| 35 MOVE STORAGE TO STORAGE                        | smove                 | SS           |
| (WITH DELTA)                                      | smovewd               | SS           |
| MOVE STORAGE TO STORAGE DIRECT                    | smoveida              | SS           |
| 40 SUBTRACT STORAGE TO STORAGE                    | ss                    | SS           |
| (WITH DELTA)                                      | sswd                  | SS           |
| SUBTRACT STORAGE TO STORAGE DIRECT                | ssda                  | SS           |
| 45 SUBTRACT STORAGE TO STORAGE FINAL              | ssf                   | SS           |
| (WITH DELTA)                                      | ssfwd                 | SS           |
| SUBTRACT STORAGE TO STORAGE FINAL DIRECT          | ssfda                 | SS           |
| 50 SUBTRACT STORAGE TO STORAGE INTERMEDIATE       | ssi                   | SS           |
| (WITH DELTA)                                      | ssiwid                | SS           |

55

Table 2 (Page 3 of 3). Storage to Storage Instructions

| <u>NAME</u>  | <u>MNE-<br/>MONIC</u> | <u>TYPE</u> |
|--|-----------------------|-------------|
| 5 SUBTRACT STORAGE TO STORAGE INTERMEDIATE<br>DIRECT   | saida                 | SS          |
| 10 SUBTRACT STORAGE TO STORAGE LOGICAL<br>(WITH DELTA) | ssl<br>sslwd          | SS<br>SS    |
| 15 SUBTRACT STORAGE TO STORAGE LOGICAL<br>DIRECT       | sslda                 | SS          |

Table 3

| Logical Instructions  |          |      |
|-----------------------|----------|------|
| NAME                  | MNEMONIC | TYPE |
| 25 AND DIRECT ADDRESS | nda      | DA   |
| AND FROM STORAGE      | n        | RS   |
| (WITH DELTA)          | nwd      | RS   |
| AND IMMEDIATE         | ni       | RI   |
| (WITH DELTA)          | nwd      | RI   |
| 30 AND REGISTER       | nr       | RR   |
| OR DIRECT ADDRESS     | oda      | DA   |
| OR FROM STORAGE       | o        | RS   |
| (WITH DELTA)          | owd      | RS   |
| OR IMMEDIATE          | oi       | RI   |
| (WITH DELTA)          | oiwd     | RI   |
| 35 OR REGISTER        | or       | RR   |
| XOR DIRECT ADDRESS    | xda      | DA   |
| XOR FROM STORAGE      | x        | RS   |
| (WITH DELTA)          | xwd      | RS   |
| 40 XOR IMMEDIATE      | xi       | RI   |
| (WITH DELTA)          | xhwd     | RI   |
| XOR REGISTER          | xr       | RR   |

Table 4 (Page 1 of 2), Shift Instructions

|    | <u>NAME</u>                                  | <u>MNE-</u><br><u>MONIC</u> | <u>TYPE</u> |
|----|--|-----------------------------|-------------|
| 5  | SCALE BINARY                                 | scale                       | SPC         |
|    | SCALE BINARY IMMEDIATE                       | scalei                      | SPC         |
|    | SCALE BINARY REGISTER                        | scalr                       | SPC         |
| 10 | SCALE HEXADECIMAL                            | scaleh                      | SPC         |
|    | SCALE HEXADECIMAL IMMEDIATE                  | scalehi                     | SPC         |
| 15 | SCALE HEXADECIMAL REGISTER                   | scalehr                     | SPC         |
|    | SHIFT LEFT ARITHMETIC BINARY                 | sla                         | SPC         |
|    | SHIFT LEFT ARITHMETIC BINARY IMMEDIATE       | slai                        | SPC         |
| 20 | SHIFT LEFT ARITHMETIC BINARY REGISTER        | slar                        | SPC         |
|    | SHIFT LEFT ARITHMETIC HEXADECIMAL            | slah                        | SPC         |
|    | SHIFT LEFT ARITHMETIC HEXADECIMAL IMMEDIATE  | slahi                       | SPC         |
| 25 | SHIFT LEFT ARITHMETIC HEXADECIMAL REGISTER   | slahr                       | SPC         |
|    | SHIFT LEFT LOGICAL BINARY                    | sl                          | SPC         |
|    | SHIFT LEFT LOGICAL BINARY IMMEDIATE          | sl                          | SPC         |
| 30 | SHIFT LEFT LOGICAL BINARY REGISTER           | slr                         | SPC         |
|    | SHIFT LEFT LOGICAL HEXADECIMAL               | slh                         | SPC         |
|    | SHIFT LEFT LOGICAL HEXADECIMAL IMMEDIATE     | slhi                        | SPC         |
| 35 | SHIFT LEFT LOGICAL HEXADECIMAL REGISTER      | slhr                        | SPC         |
|    | SHIFT RIGHT ARITHMETIC BINARY                | sra                         | SPC         |
|    | SHIFT RIGHT ARITHMETIC BINARY IMMEDIATE      | srai                        | SPC         |
|    | SHIFT RIGHT ARITHMETIC BINARY REGISTER       | srar                        | SPC         |
| 40 | SHIFT RIGHT ARITHMETIC HEXADECIMAL           | srah                        | SPC         |
|    | SHIFT RIGHT ARITHMETIC HEXADECIMAL IMMEDIATE | srahi                       | SPC         |
| 45 | SHIFT RIGHT ARITHMETIC HEXADECIMAL REGISTER  | srahr                       | SPC         |
|    | SHIFT RIGHT LOGICAL BINARY                   | srl                         | SPC         |
| 50 | SHIFT RIGHT LOGICAL BINARY IMMEDIATE         | srl                         | SPC         |

55

Table 4 (Page 2 of 2). Shift Instructions

|    | <u>NAME</u>                               | <u>MNE-</u><br><u>MONIC</u> | <u>TYPE</u> |
|----|---|-----------------------------|-------------|
| 5  | SHIFT RIGHT LOGICAL BINARY REGISTER       | srlr                        | SPC         |
|    | SHIFT RIGHT LOGICAL HEXADECIMAL           | srlh                        | SPC         |
|    | SHIFT RIGHT LOGICAL HEXADECIMAL IMMEDIATE | srlhi                       | SPC         |
| 10 | SHIFT RIGHT LOGICAL HEXADECIMAL REGISTER  | srlhr                       | SPC         |

Table 5 (Page 1 of 2). Branch Instructions

|    | <u>NAME</u>               | <u>MNE-</u><br><u>MONIC</u> | <u>TYPE</u> |
|----|---------------------------|-----------------------------|-------------|
| 20 | BRANCH                    | b                           | RS          |
|    | (WITH DELTA)              | bwd                         | RS          |
|    | BRANCH DIRECT             | bda                         | DA          |
| 25 | BRANCH IMMEDIATE          | bi                          | RI          |
|    | (WITH DELTA)              | biwd                        | RI          |
|    | BRANCH REGISTER           | br                          | RS          |
|    | BRANCH AND LINK           | bal                         | RS          |
| 30 | BRANCH AND LINK DIRECT    | balda                       | DA          |
|    | BRANCH AND LINK IMMEDIATE | ball                        | RI          |
|    | (WITH DELTA)              | ballwd                      | RI          |
| 35 | BRANCH AND LINK REGISTER  | balr                        | RS          |
|    | BRANCH BACKWARD           | bb                          | RS          |
|    | (WITH DELTA)              | bbwd                        | RS          |
| 40 | BRANCH BACKWARD DIRECT    | bbda                        | DA          |
|    | BRANCH BACKWARD IMMEDIATE | bbi                         | RI          |
|    | (WITH DELTA)              | bbiwd                       | RI          |
| 45 | BRANCH BACKWARD REGISTER  | bbr                         | RS          |
|    | BRANCH FORWARD            | bf                          | RS          |
|    | (WITH DELTA)              | bfwd                        | RS          |
| 50 | BRANCH FORWARD DIRECT     | bfda                        | DA          |



Table 5 (Page 2 of 2). Branch Instructions

| <u>NAME</u>                         | <u>MNE-<br/>MONIC</u> | <u>TYPME</u> |
|-------------------------------------|-----------------------|--------------|
| 5     BRANCH FORWARD IMMEDIATE      | bfi                   | RI           |
| (WITH DELTA)                        | bfiwd                 | RI           |
| BRANCH FORWARD REGISTER             | bfr                   | RS           |
| 10    BRANCH ON CONDITION           | bc                    | RS           |
| (WITH DELTA)                        | bcwd                  | RS           |
| BRANCH ON CONDITION DIRECT          | bcda                  | RS           |
| 15    BRANCH ON CONDITION IMMEDIATE | bci                   | RI           |
| (WITH DELTA)                        | bciwd                 | RI           |
| BRANCH ON CONDITION REGISTER        | bcr                   | RS           |
| 20    BRANCH RELATIVE               | brel                  | RI           |
| (WITH DELTA)                        | brelwd                | RS           |
| NULL OPERATION                      | noop                  | RR           |

Table 6

| Status Switching Instructions |          |       |
|-------------------------------|----------|-------|
| NAME                          | MNEMONIC | TYPME |
| RETURN                        | ret      | SPC   |

Table 7

| Input/Output Instructions |          |       |
|---------------------------|----------|-------|
| NAME                      | MNEMONIC | TYPME |
| IN                        | IN       | SPC   |
| OUT                       | OUT      | SPC   |
| INTERNAL DIOR/DIOW        | INTR     | SPC   |

### SOME SUMMARY FEATURES

#### The APAP Machine in Perspective

50     We have described in accordance with our invention could be thought of in its more detailed aspects to be positioned in the technology somewhere between the CM-1 and N-cube. Like our APAP, the CM-1 uses a point design for the processing element and combines processing elements with memory on the basic chip. The CM-1, however uses a 1 bit wide serial processor, while the APAP series will use a 16 bit wide processor. The CM series of machines started with 4K bits of memory per processor and has grown to 8 or 16K bits versus the 32K by 16 bits we have provided for the first APAP chip. The CM-1 and its follow-ons are strictly SIMD machines while the CM-5 is a hybrid. Instead of this, our APAP will effectively use MIMD operating modes in conjunction with SIMD modes when useful. While our parallel 16 bit wide PMEs might

be viewed as a step toward the N-cube, this step is not warranted. The APAP does not separate memory and routing from the processing element as does the N-cube kind of machine. Also, the APAP provides for up to 32K 16 bit PMEs while the N-cube only provides for 4K 32 bit processors.

Even with the superficial similarities presented above, the APAP concept completely differs from the CM and N-cube series by:

1. The modified hypercube incorporated in our APAP is a new invention providing a significant packaging and addressing advantage when compared with hypercube topologies. For instance, consider that the 32K PME APAP in its first preferred embodiment has a network diameter of 19 logical steps and, with transparency, this can be reduced to an effective 16 logical steps. Further, by comparison, if a pure hypercube were used, and if all PMEs were sending data through an 8 step path, then on average 2 of every 8 PMEs would be active while the remainder would be delayed due to blockage.

Alternatively, consider the 64K hypercube that would be needed if CM-1 was a pure hypercube. In that case, each PME would require ports to 16 other PMEs, and data could be routed between the two farthest separated PMEs in 15 logical steps. If all PMEs tried to transfer an average distance of 7 steps, the 2 of every 7 would be active. However, CM-1 does not utilize a 16d hypercube. It interconnects the 18 nodes on a chip with a NEWS network; then it provides one router function within the chip. The 4096 routers are connected into a 12d hypercube. With no collisions the hybrid still has a logical diameter of 15, but since 16 PMEs could be contending for the link its effective diameter is much greater. That is, with 8 step moves only 2 of 16 PMEs could be active, which means that 8 complete cycles rather than 4 cycles are needed to complete all data moves.

The N-cube actually utilizes a pure hypercube, but currently only provides for a 4096 PMEs and thus, utilizes a 12d (13d for 8192 PMEs) hypercube. For the N-cube to grow to 16K processors, at which point it would have the same processing data width as the APAP, it would have to add four times as much hardware and would have to increase the connection ports to each PME router by 25%. Although no hard data exists to support this conclusion, it would appear that the N-cube architecture runs out of connector pins prior to reaching a 16K PME machine.

2. The completely integrated and distributed nature of major tasks within the APAP machine is a decided advantage. As was noted for the CM and N-cube series of machines, each had to have separate units for message routing as well as separate units for floating point coprocessors. The APAP system combines the integer, floating point processing, message routing and I/O control into the single point design PME. That design is then replicated 8 times on a chip, and the chip is then replicated 4K times to produce the array. This provides several advantages:

- a. Using one chip means maximum size production runs and minimal system factor costs.
  - b. Regular architecture produces the most effective programming systems.
  - c. Almost all chip pins can be dedicated to the generic problem of interprocessor communication, maximizing the inter-chip I/O bandwidth which tends to be a important limiting factor in MPP designs.
3. The APAP has the unique design ability to take advantage of chip technology gains and capital investment in custom chip designs.

Consider the question of floating point performance. It is anticipated that APAP PME performance on DAXPY will be about 125 cycles per flop. In contrast, the '387 Coprocessor would be about 14 cycles while the Weitek Coprocessor in the CM-1 would be about 6 cycles. However, in the CM case there is only one floating point unit for every 16 PMEs while in the N-cube case there is probably one '387 type chip associated with each of the '386 processors. Our APAP has 16 times as many PMEs and therefore can almost completely make up for the single unit performance delta.

More significantly, the 8 APAP PMEs within a chip are constructed from 50K gates currently available in the technology. As memory macros shrink and the number of gates available to the logic increases, spending that increase on enhanced floating point normalization should permit APAP floating point performance to far exceed the other units. Alternatively, effort could be spent to generate a PME or PME subsection design using custom design approaches, enhancing total performance while in no way affecting any SAV developed for the machine.

We believe our design for our APAP has characteristics poised to take advantage of the future process technology growth. In contrast, the nearest similar machines CM-x and N-cube which employ a system like that described in FIGURE 1 seem well poised to take advantage of yesterday's technology which we feel is dead ended.

An advantage of the APAP concept is the ability to use DASD associated with groups of PMEs. This APAP capability, as well as the ability to connect displays and auxiliary storage, is a by-product of picking MC bus structures as the interface to the external I/O ports of the PME Array. Thus, APAP systems will be configurable and can include card mounted hard drives selected from one of the set of units that are

compatible with PS/2 or RISC/6000 units. Further, that capability should be available without designing any additional part number modules although it does require utilizing more replications of the backpanel and base enclosure than does the APAP.

This brief perspective is not intended to be limiting, but rather is intended to cause those skilled in the art to review the foregoing description and examine how the many inventions we have described which may be used to move the art of massively parallel systems ahead to a time when programming is no longer a significant problem and the costs of such systems are much lower. Our kind of system can be made available, not only to the few, but to many as it could be made at a cost within the reach of commercial department level procurements.

While we have described our preferred embodiments of our invention, it will be understood that those skilled in the art, both now and in the future, upon the understanding of these discussions will make various improvements and enhancements thereto which fall within the scope of the claims which follow. Those claims should be construed to maintain the proper protection for the invention first disclosed.

#### 15 Claims

1. A computer system comprising, a plurality of multi-processor memory elements, each having communication paths, processor and memory, and wherein a programmable router is provided for routing data and control information from one multi-processor memory element to another multi-processor memory element and between nodes of the computer system.
2. A computer system according to claim 1 wherein each multi-processor memory element (PME) has 2n processors, and communication paths which minimize delays due to chip crossings.
3. A computer system according to claim 1 wherein each multi-processor memory element (PME) has a processor, memory and routers within a single chip and internal and external communication paths which minimize delays due to chip crossings, each processor memory element having means for fixed and floating point processing, routing and I/O control.
4. A computer system according to claim 1 further comprising within a processor memory element: a native instruction set, means for providing an expandable multiply function, a programmer router for routing information alternatively left/right, NEWS matrix, NEWS/shut-down, hypercube, and wherein said programmable router is employs a hardwired distributed router provided by each processing memory element.
5. A computer system according to claim 1 organized as a massively parallel machine with nodes interconnected as a n dimensional network cluster with parallel communication paths between processor memory elements along said internal and external communication paths providing a processing array, and wherein processing memory elements of an array have a transparent mode utilized when routing data between processing memory elements within a chip set of processing memory elements for permitting reduction of the effective network diameter of a network of nodes.
6. A computer system according to claim 1 wherein a node of a processor array has multiple single processor elements made up of 32K 16-bit words with a 16-bit processor for a network node of eight processors with their associated memory with their fully distributed I/O routers and signal I/O ports.
7. A computer system according to claim 1 wherein a node of a processor array has multiple single processor elements made up of 32K 16-bit words with a 16-bit processor for a network node of eight processors with their associated memory with their fully distributed I/O routers and signal I/O ports combined as groups groups of node clusters organized as a 2d modified hypercube.
8. A computer system according to claim 1 wherein a node of a processor array has multiple single processor elements made up of 32K 16-bit words with a 16-bit processor for a network node of eight processors with their associated memory with their fully distributed I/O routers and signal I/O ports combined as groups groups of node clusters organized as a 2d modified hypercube, with up to 64 clusters integrated in a network of node clusters to form are integrated to form a 4d modified hypercube of up to 32,768 processing memory elements.

9. A computer system according to claim 1 wherein a node processing memory element has internal data flows using high speed hard registers to feed distributed ALU and and I/O router registers and logic for all operations.
- 5 10. A computer system according to claim 1 wherein a node processing memory element in has an I/O port for off chip byte wide communication, and has input ports that are connected such that data may be routed from input to memory, or from an input address register to an output register via a direct parallel data path.
- 20 11. A computer system according to claim 1 wherein a node has multiple processor memory elements and is connected to other nodes in a cluster network with data routing distributed between hardware and software, with software controlling most of the task sequencing function.
- 25 12. A computer system according to claim 1 wherein a node has multiple processor memory elements and is connected to other nodes in a cluster network with data routing distributed between hardware and software, with hardware provided for performing inner loop transfers and minimizing overhead on the outer loops of the node.
- 30 13. A computer system according to claim 1 wherein a node has multiple processor memory elements and is connected to other nodes in a cluster network with I/O programs at dedicated interrupt levels for managing the network.
- 35 14. A computer system according to claim 1 wherein a node has multiple processor memory elements and is connected to other nodes in a cluster network with I/O programs at dedicated interrupt levels for managing the network, each processor memory element having interrupt registers and dedicating four interrupt levels to receiving data from four neighbors, a buffer provided at each level by loading registers at the level, and having in and return instruction pairs using a buffer address and transfer count to enable the processor memory element to accept words from an input bus and to store them to the buffer.
- 40 15. A multi-processor memory system, comprising: a plurality of multi-processor memory elements, each multi-processor memory element (PME) having 2n processors, memory and routers within a single chip and internal and external communication paths which minimize delays due to chip crossings.

**FIG.1A**  
Prior Art

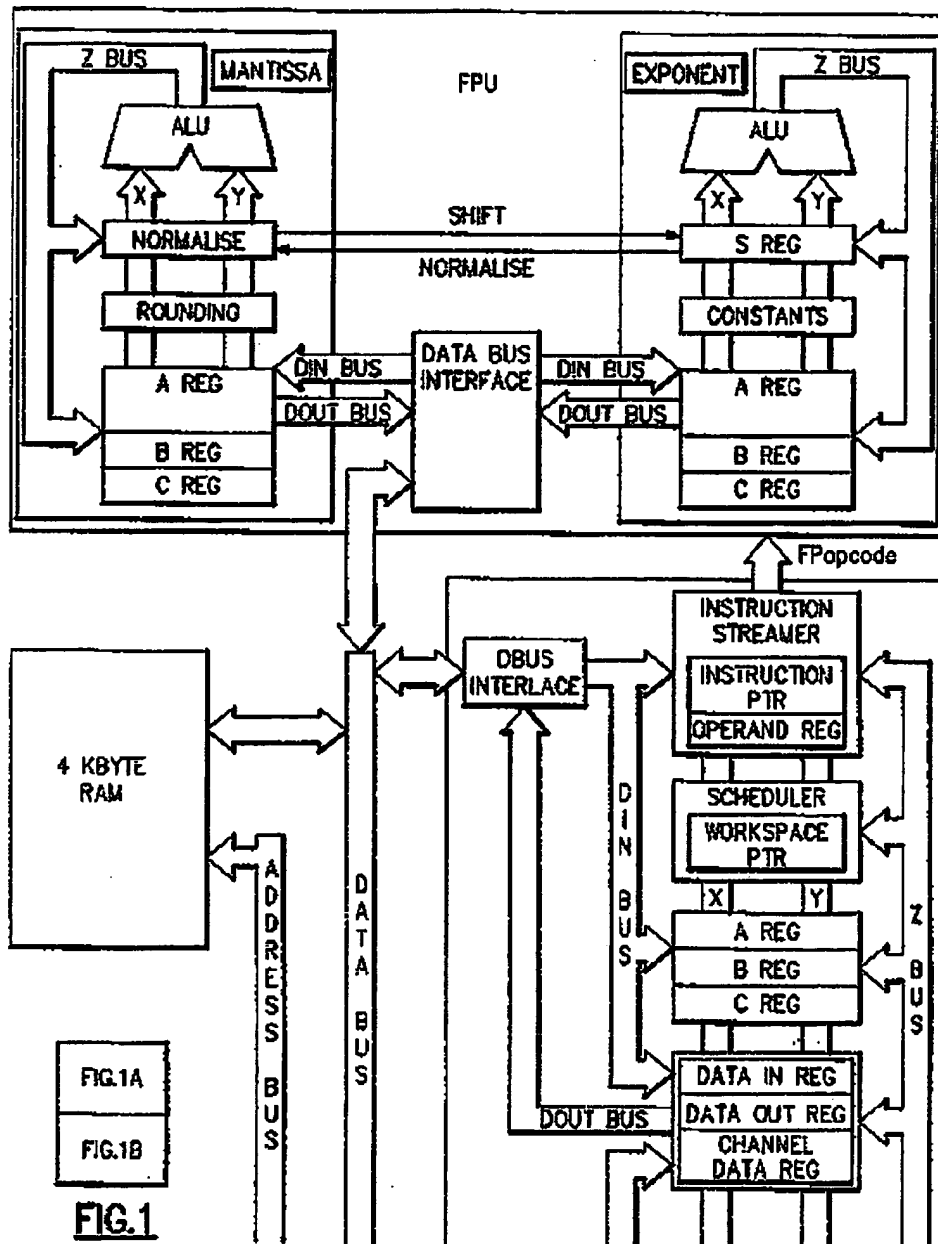
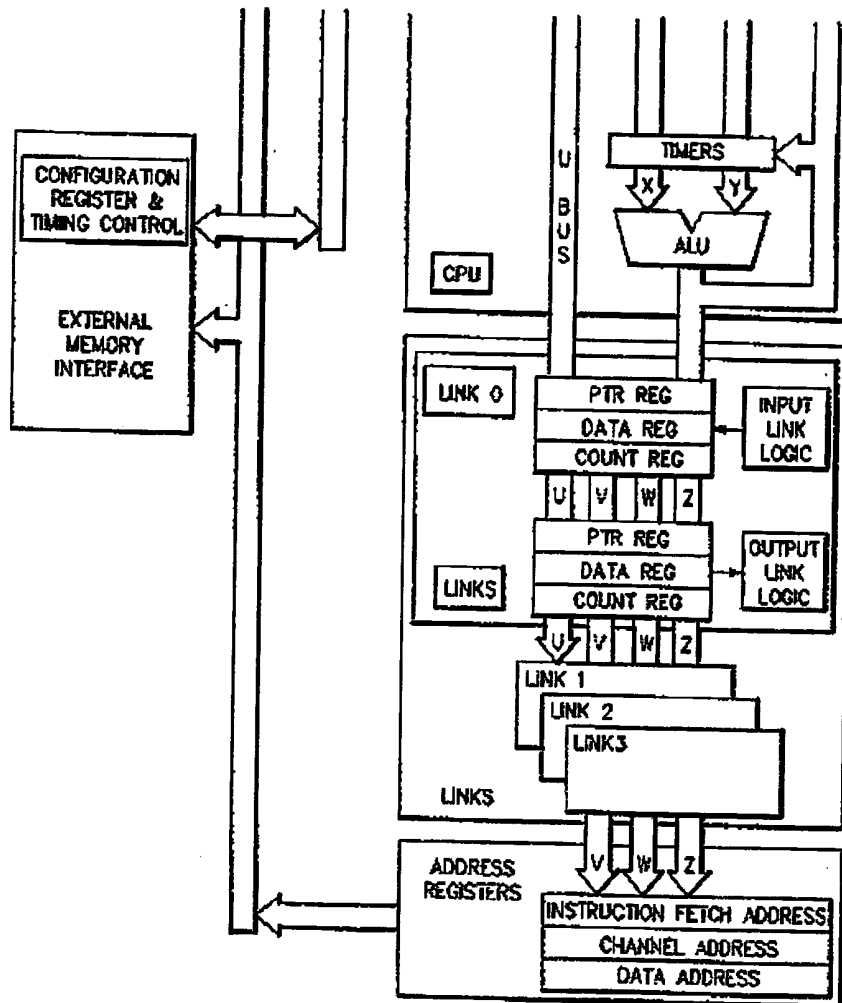


FIG.1A

FIG.1B

**FIG.1**



**FIG.1B**  
Prior Art

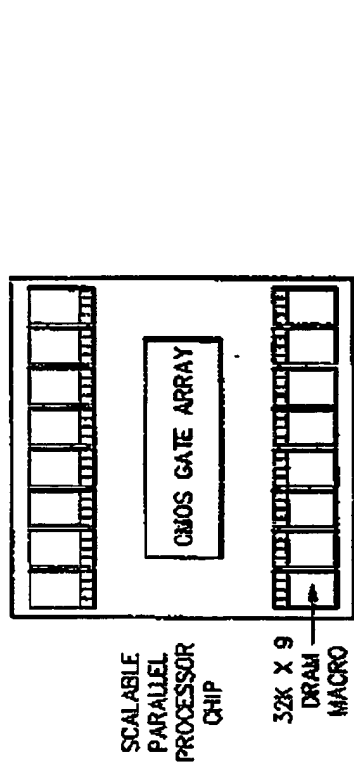


FIG. 2

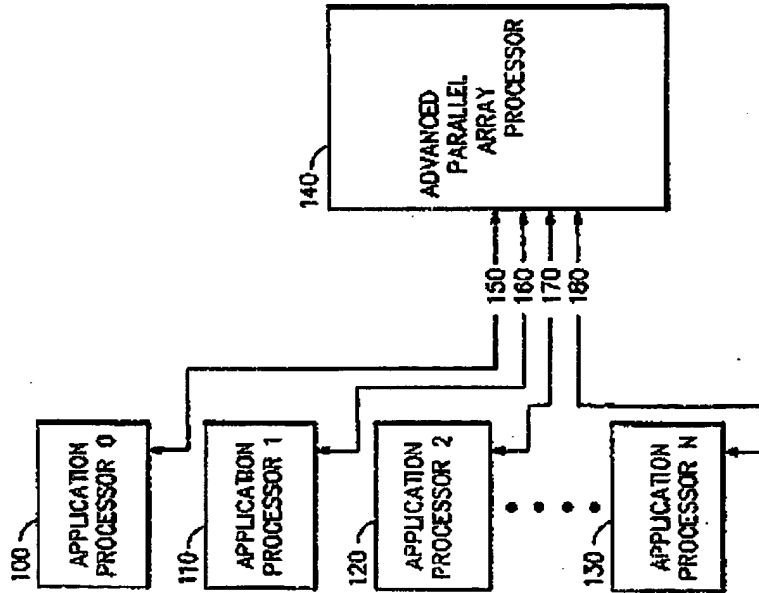


FIG. 4

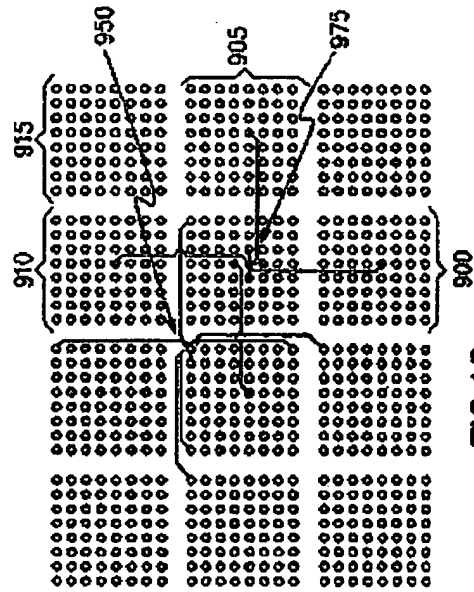
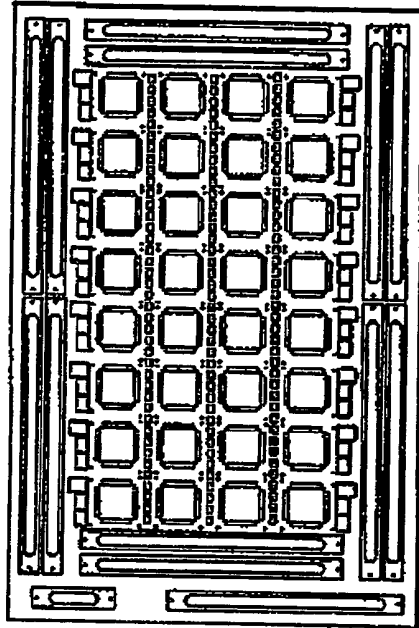


FIG. 16



**FIG. 3**



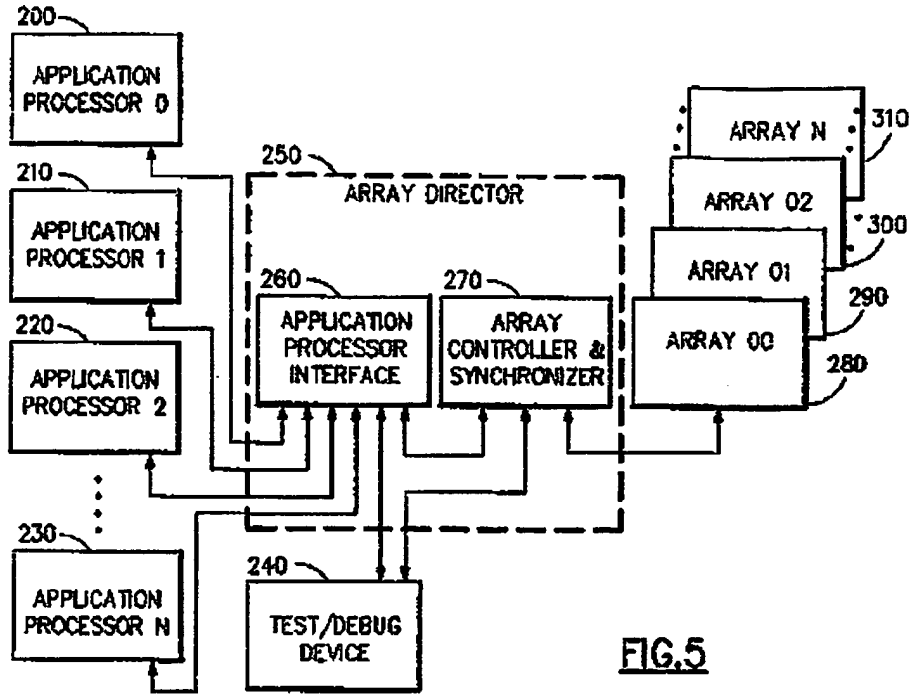


FIG. 5

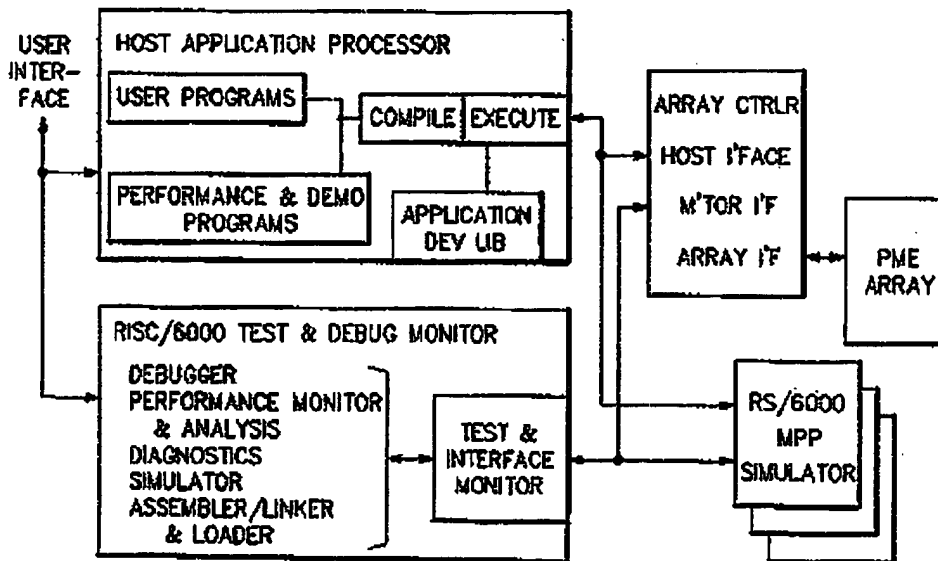


FIG. 19

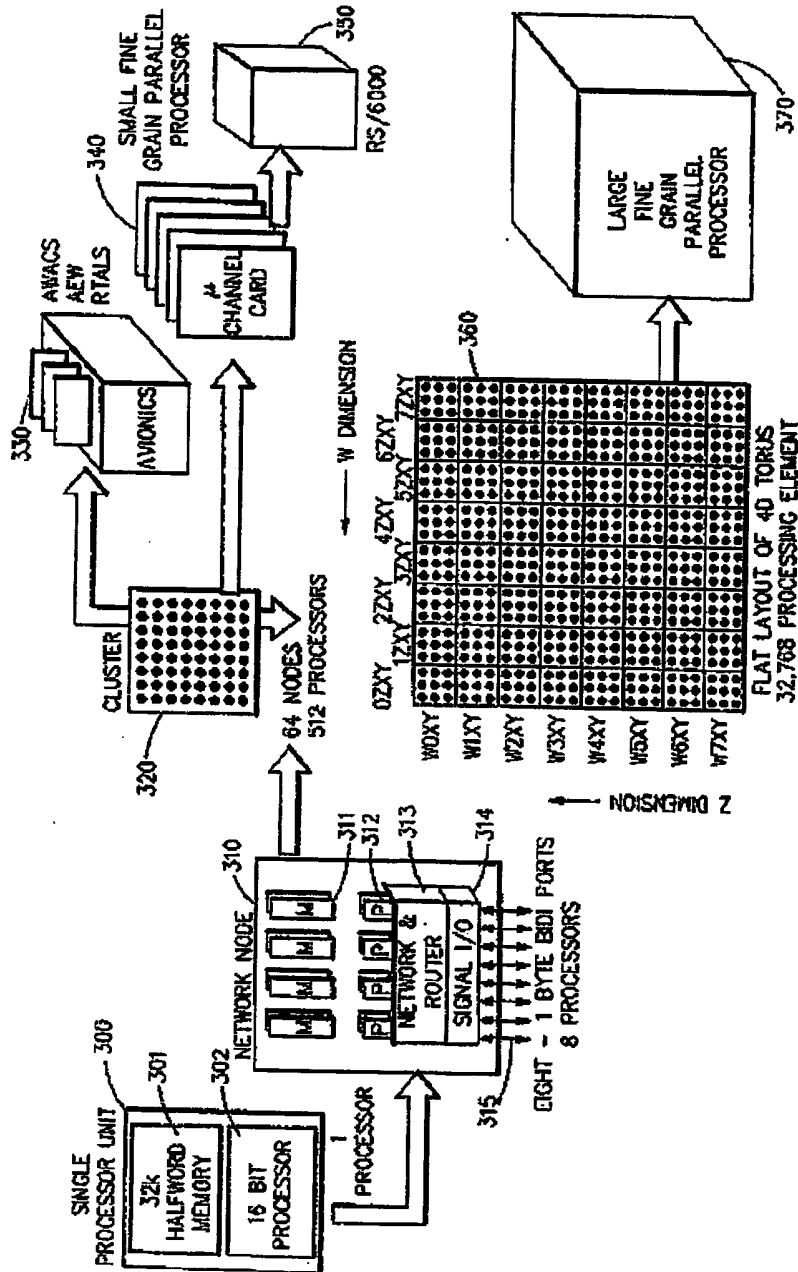


FIG.6

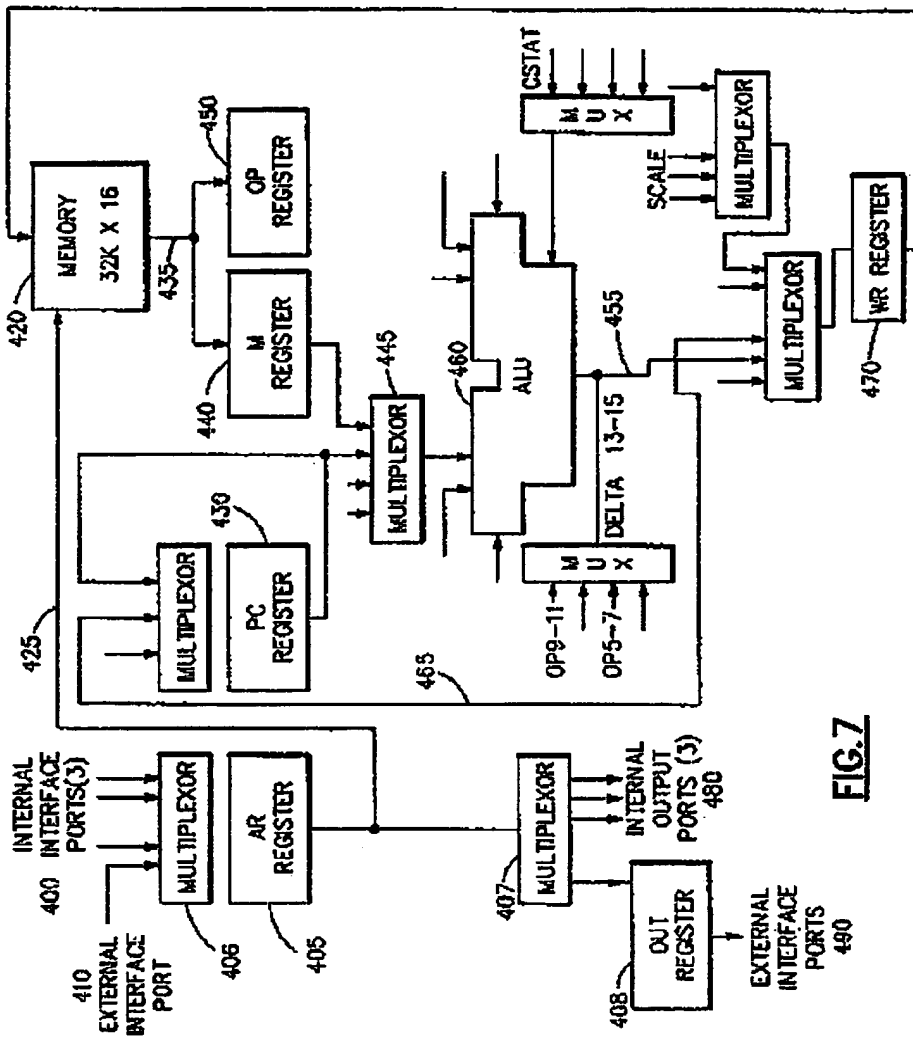


FIG. 7

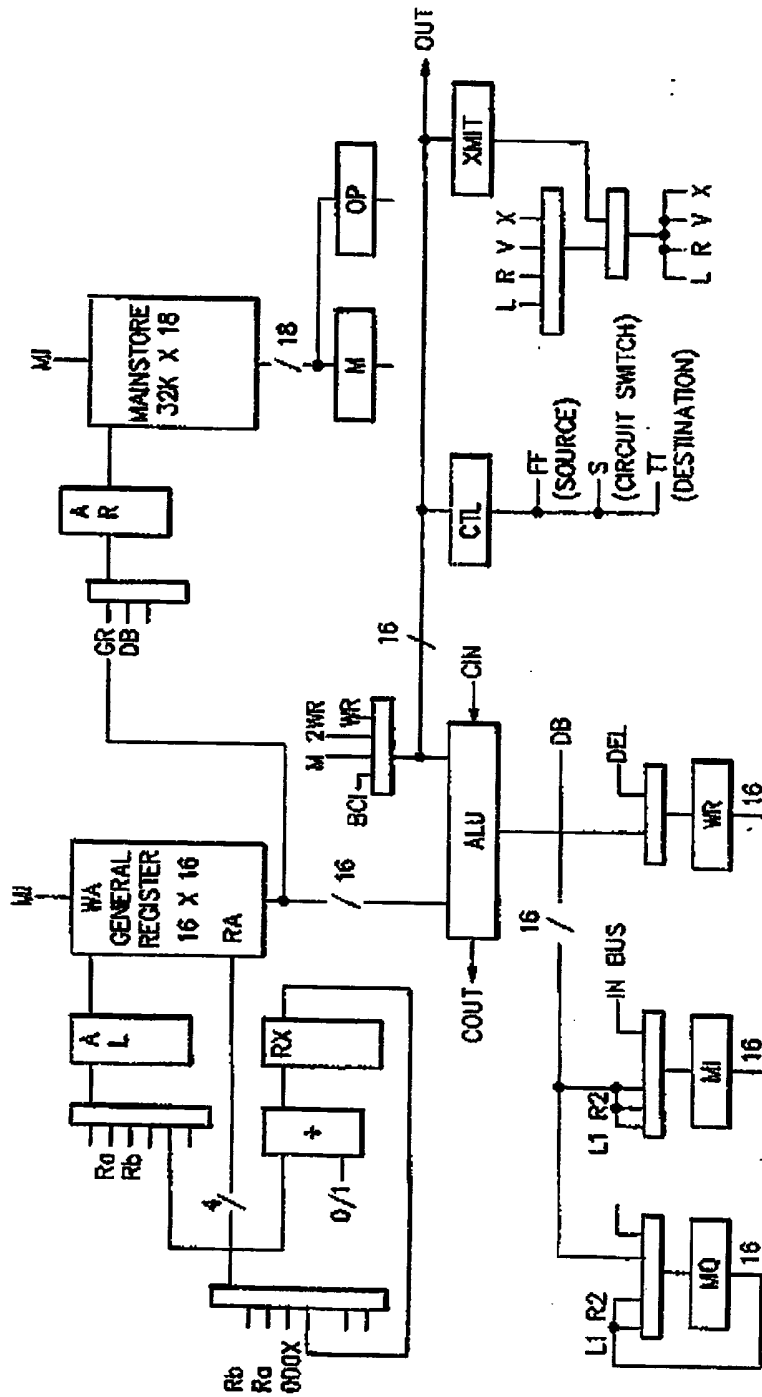
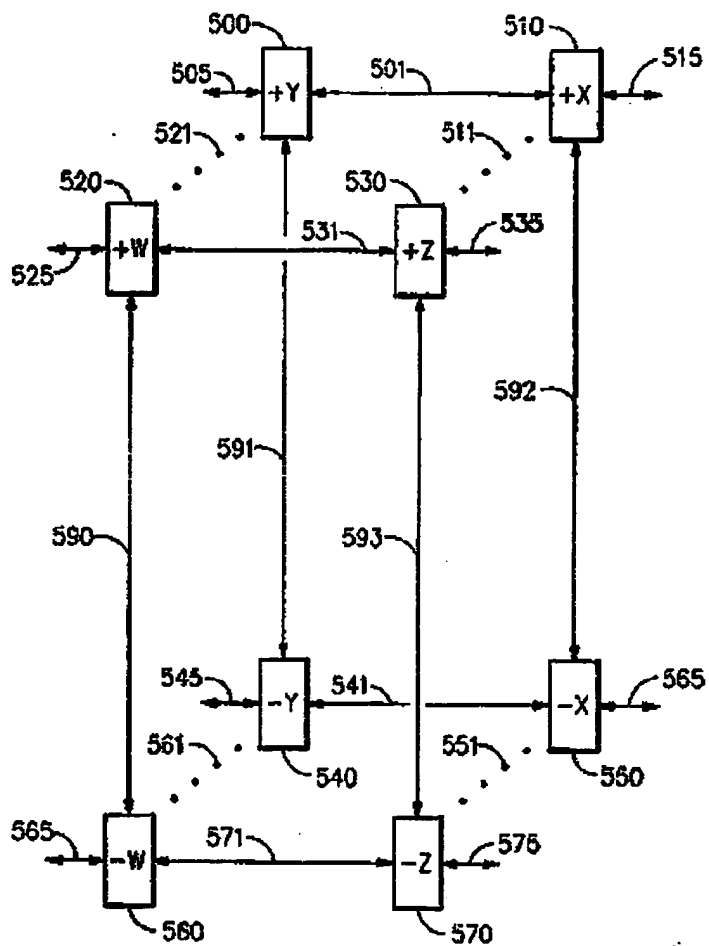


FIG. 8



**FIG. 9**

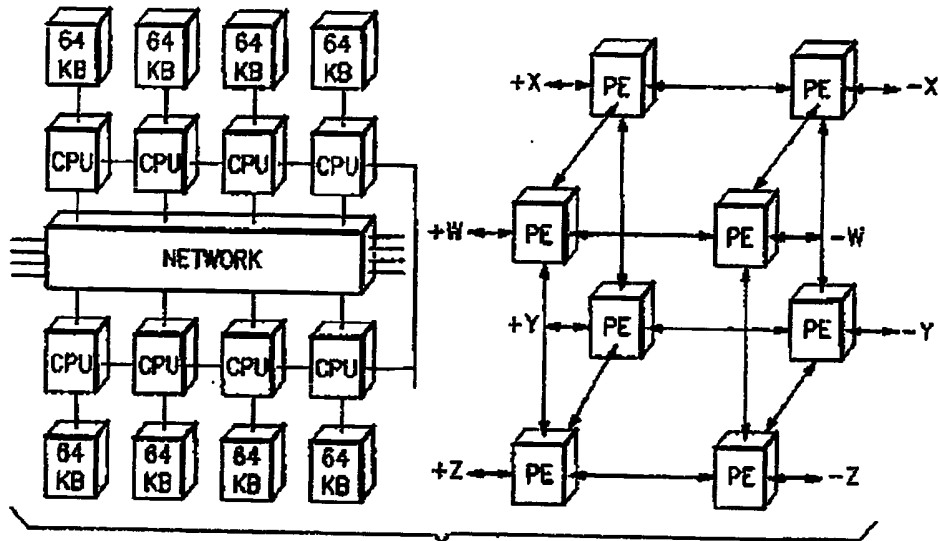


FIG. 10

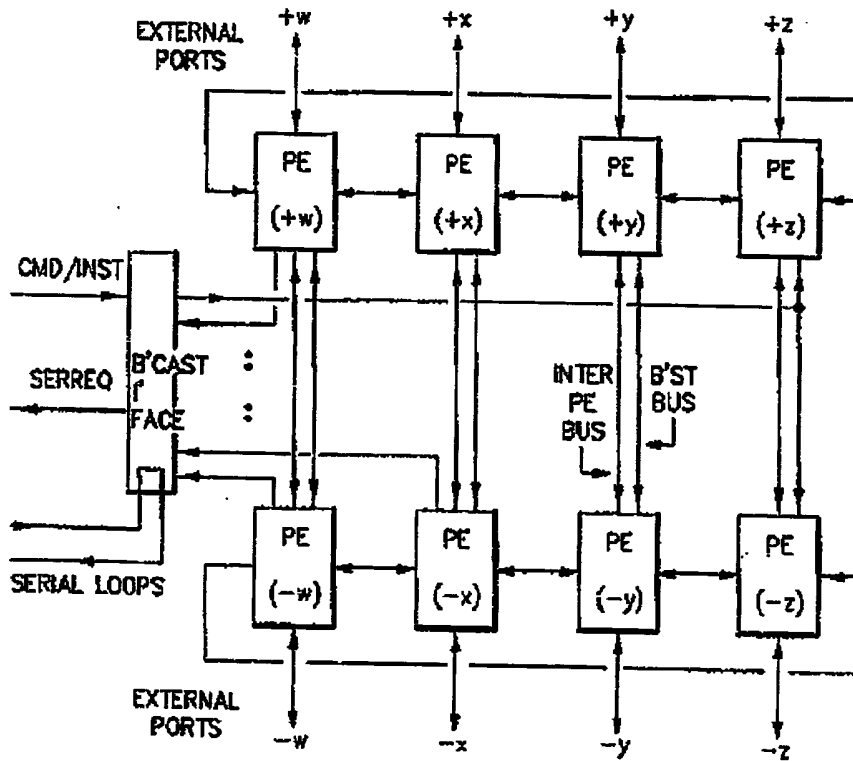


FIG. 11

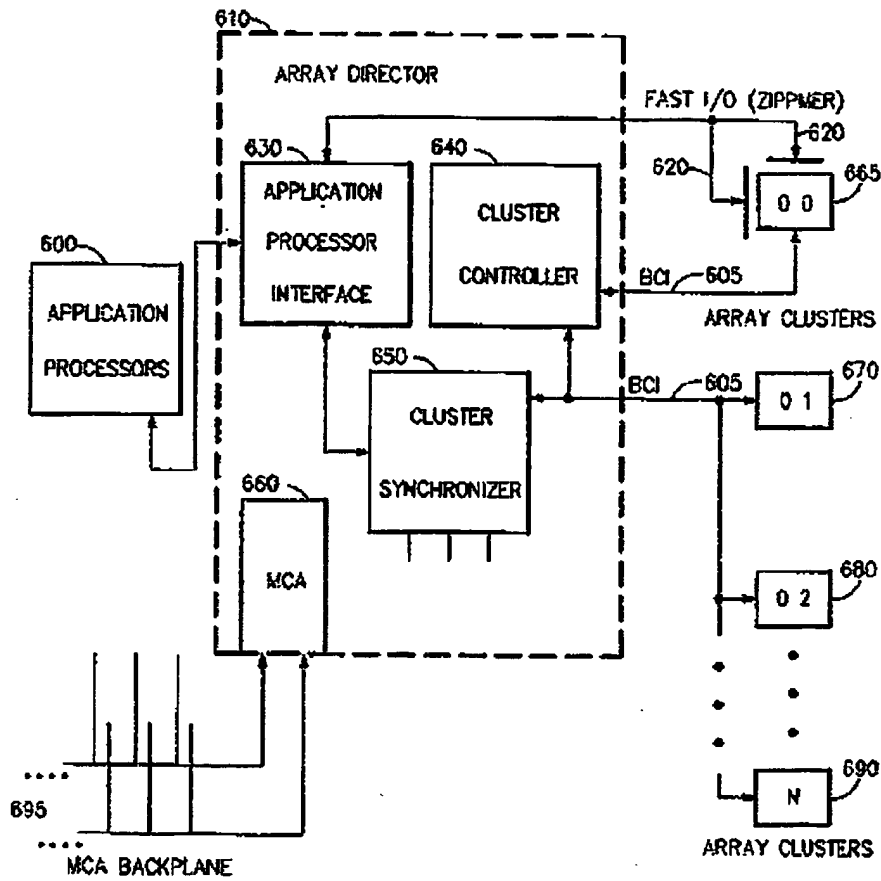


FIG.12

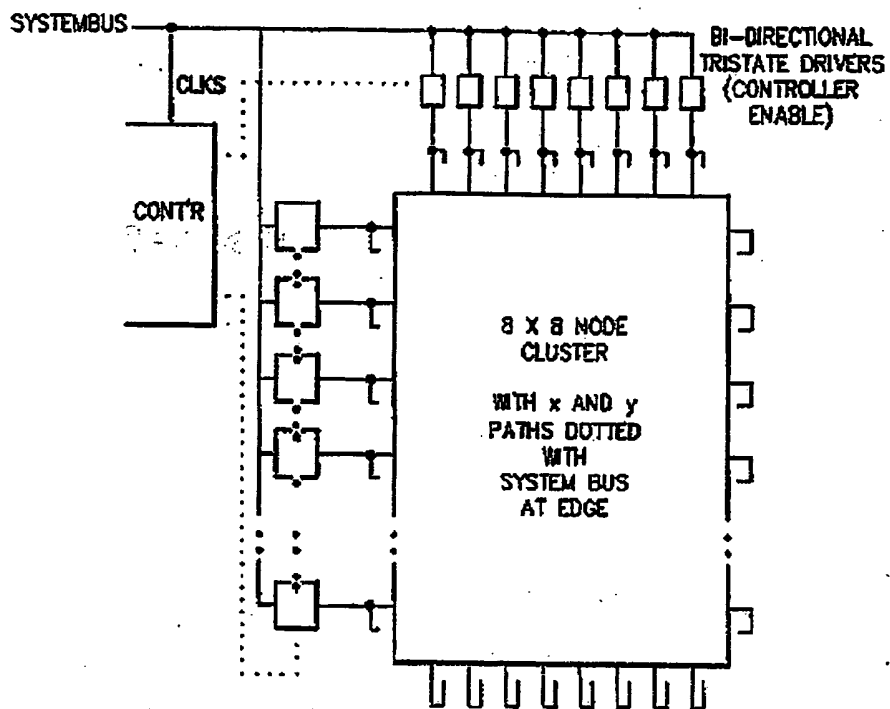


FIG.13A

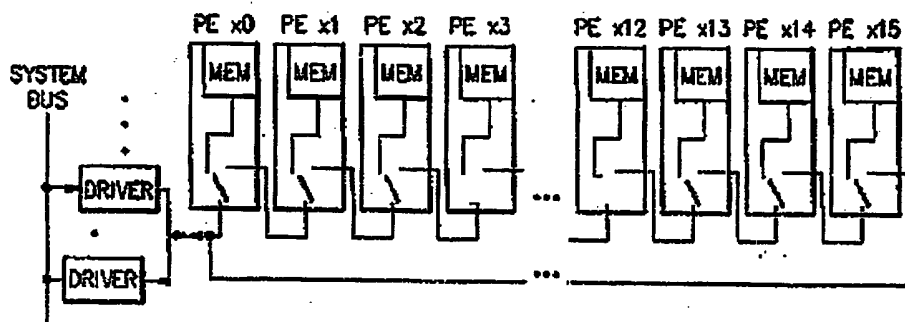


FIG.13B



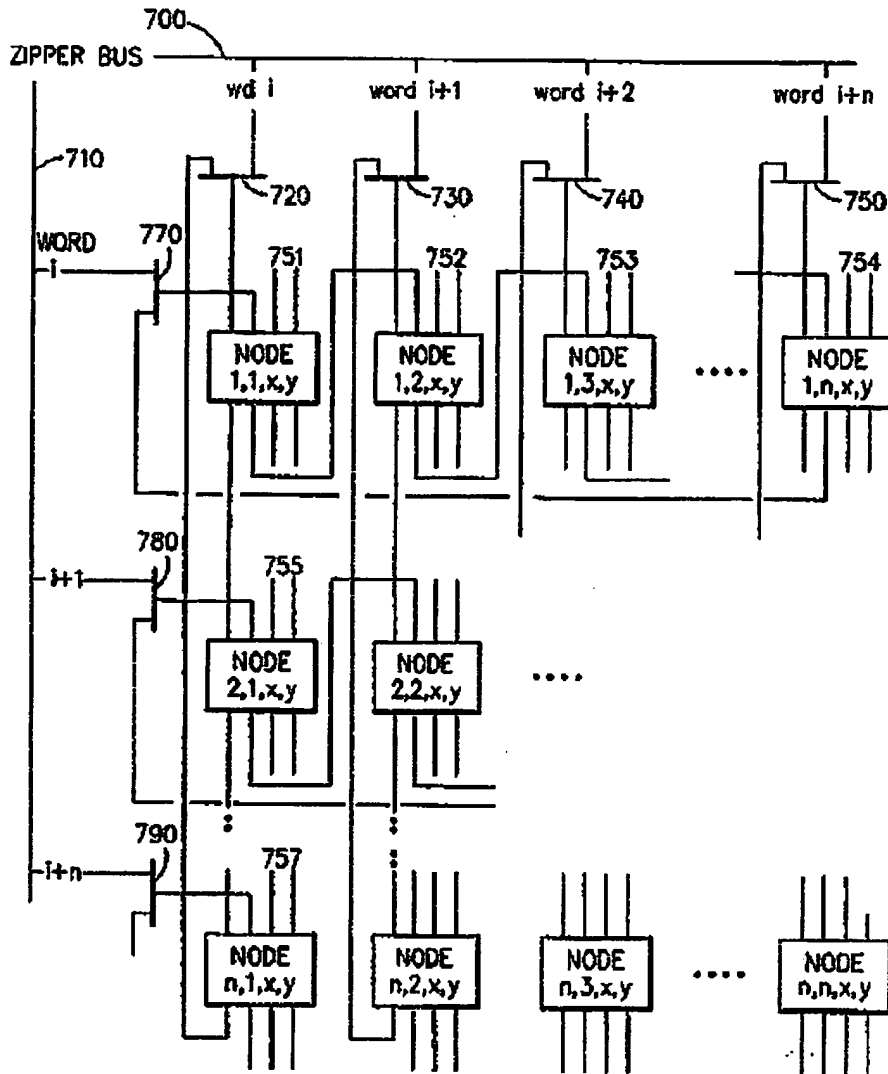


FIG.14

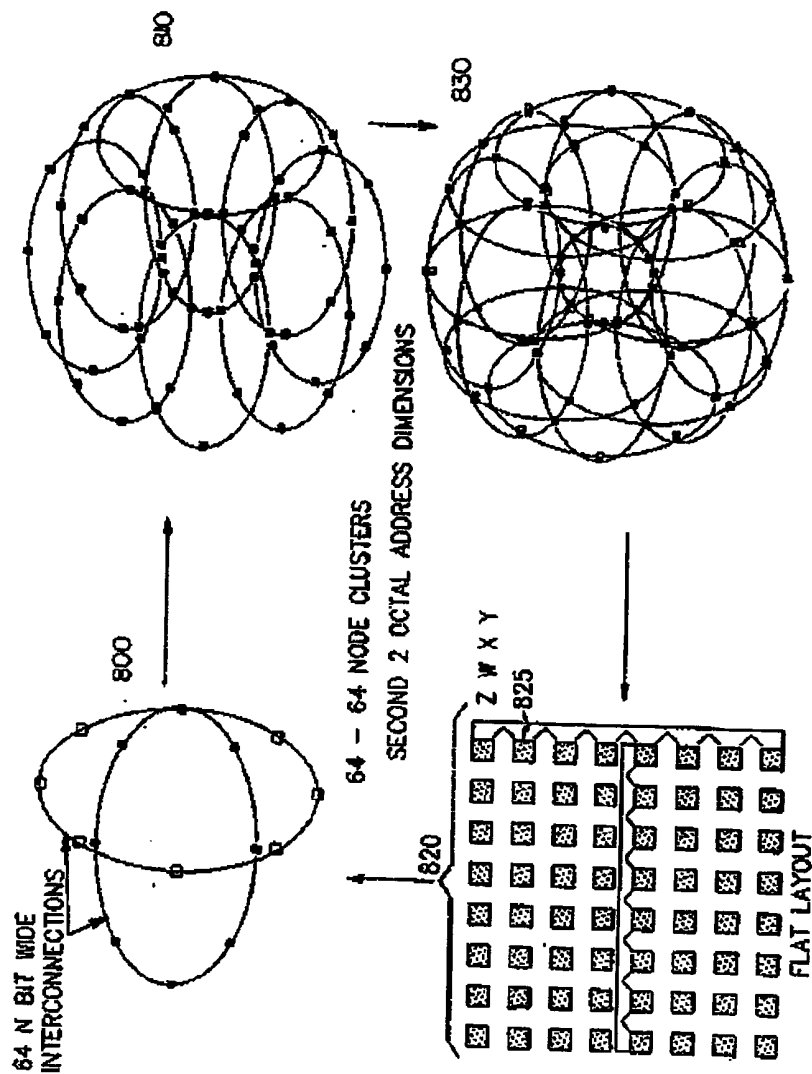
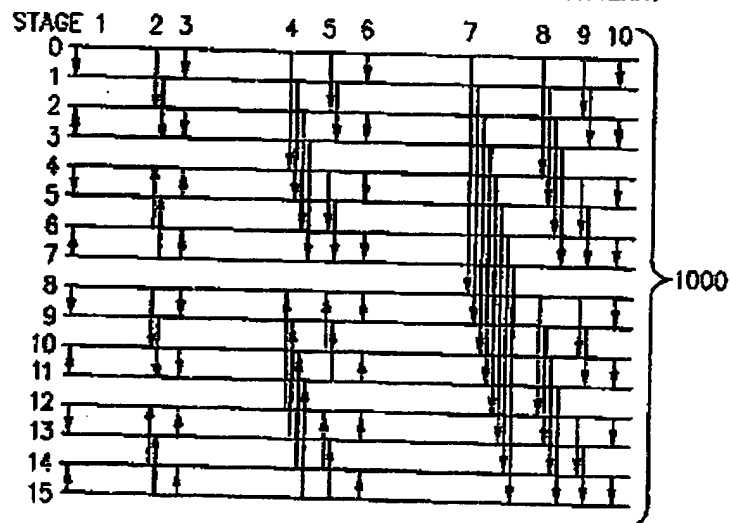


FIG.15

HOW WOULD A 16 ELEMENT SORT REPEAT THE PATTERN?



1100 { FOR SORTING  $n$  DATA ELEMENTS ( $n \in \{2^i : i \in \mathbb{N}, 2^i \leq \# \text{ OF PE'S}\}$ )

do  $i = 0$  to  $(\log_2 n) - 1$  do  $j = 0$  to  $i$

if  $(\text{PE\#} / 2^i - j) \% 2 = 0$

then  $\text{TARGET} = \text{PE\#} + 2^i - j$  else  $\text{TARGET} = \text{PE\#} - 2^i - j$

send DATA to TARGET

receive data store in TEMP (if data is not available - wait)

if  $(2((\frac{\text{PE\#}}{2^{j+1}}) \% 2) + ((\frac{\text{PE\#}}{2^{j+1}}) \% 2) + 1) \% 2 = 0$

then if  $\text{TEMP} < \text{DATA}$  then  $\text{DATA} = \text{TEMP}$  else NOP

then if  $\text{TEMP} > \text{DATA}$  then  $\text{DATA} = \text{TEMP}$  else NOP

end both do's

FIG.17

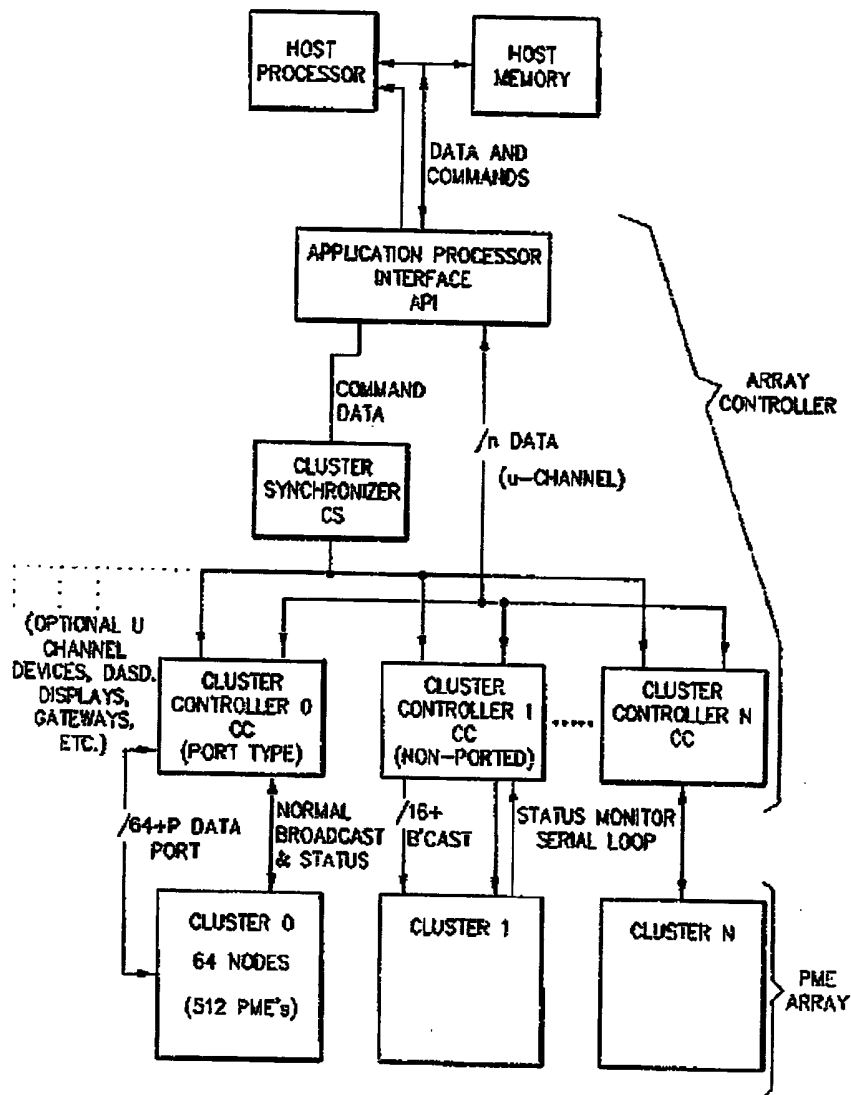


FIG.18

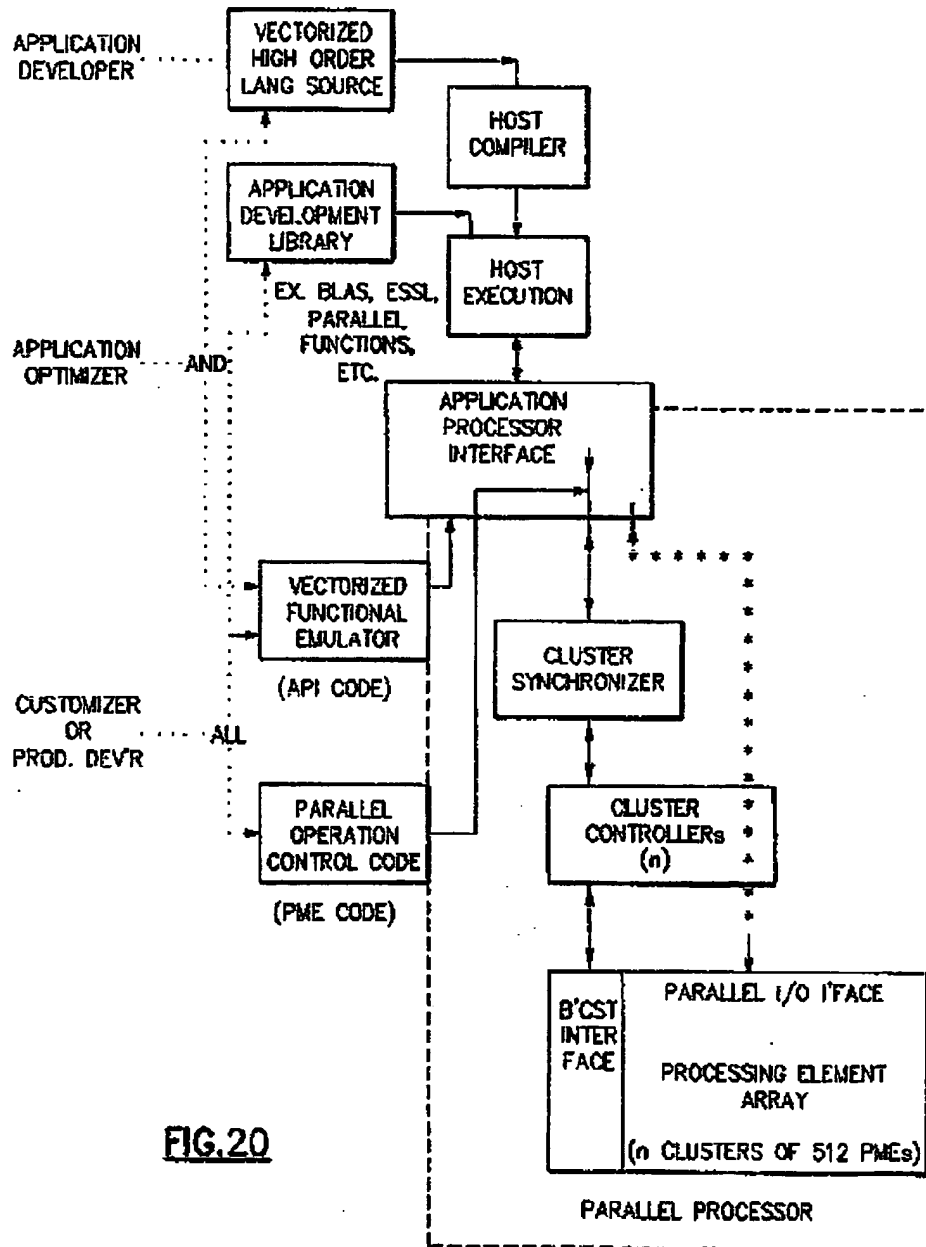


FIG.20

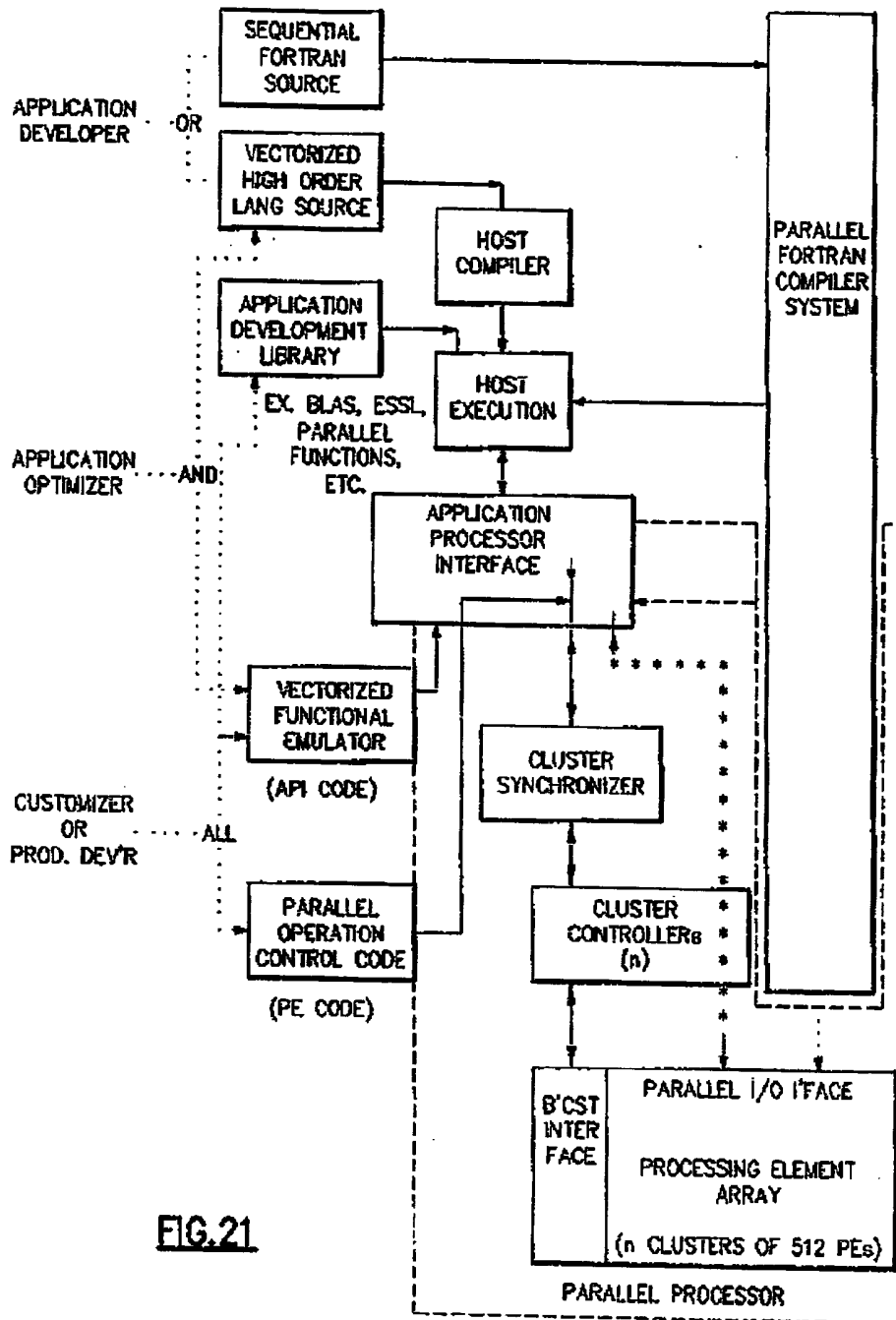


FIG.21

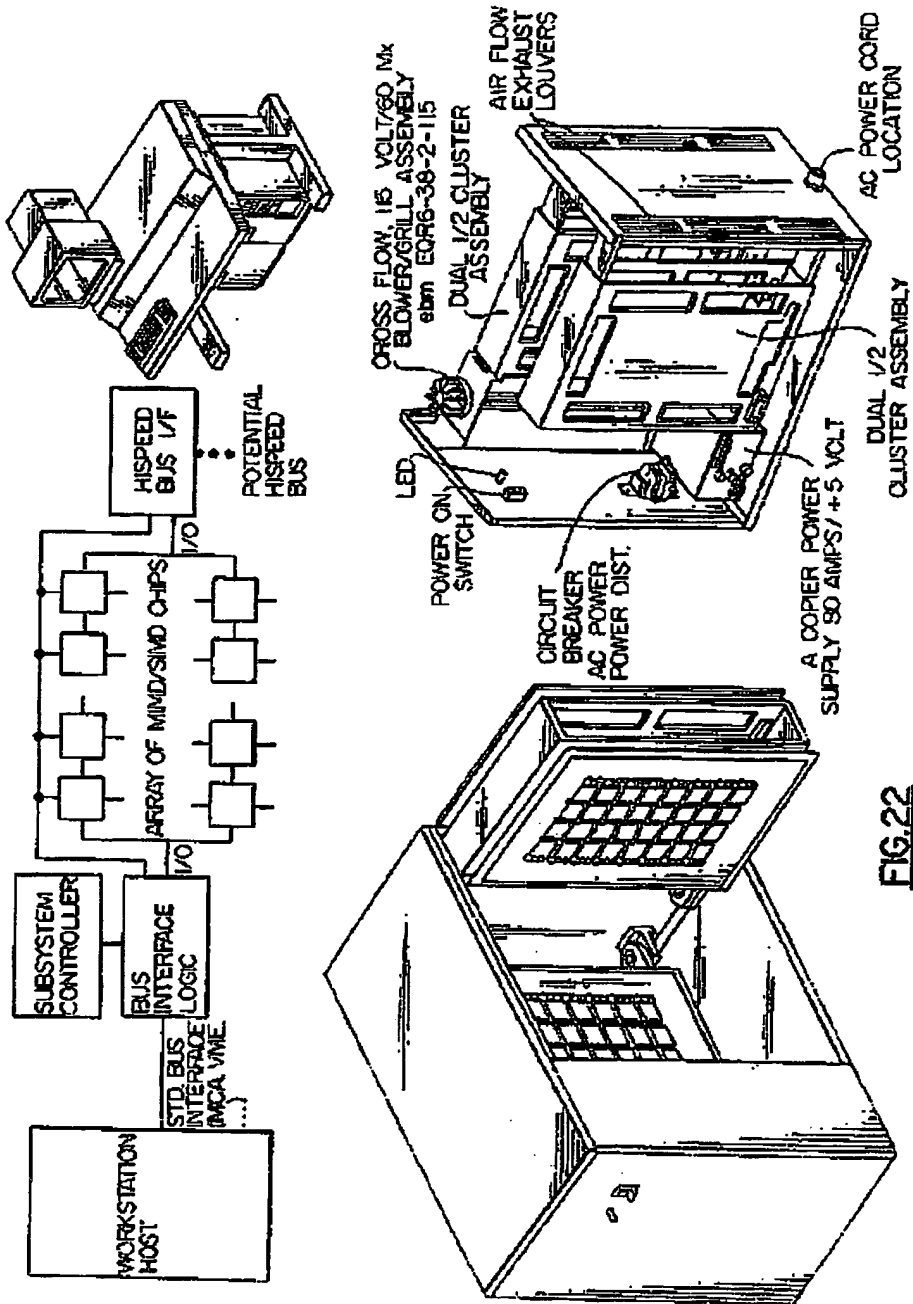


FIG. 22

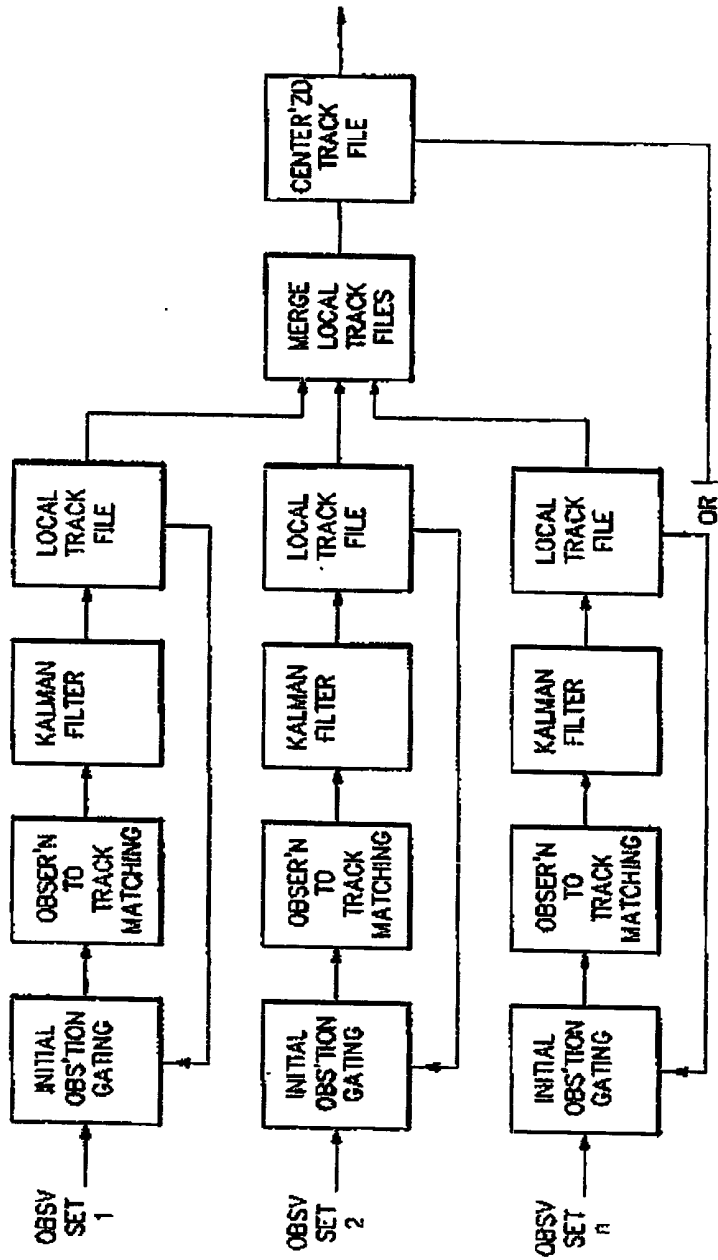


FIG. 23



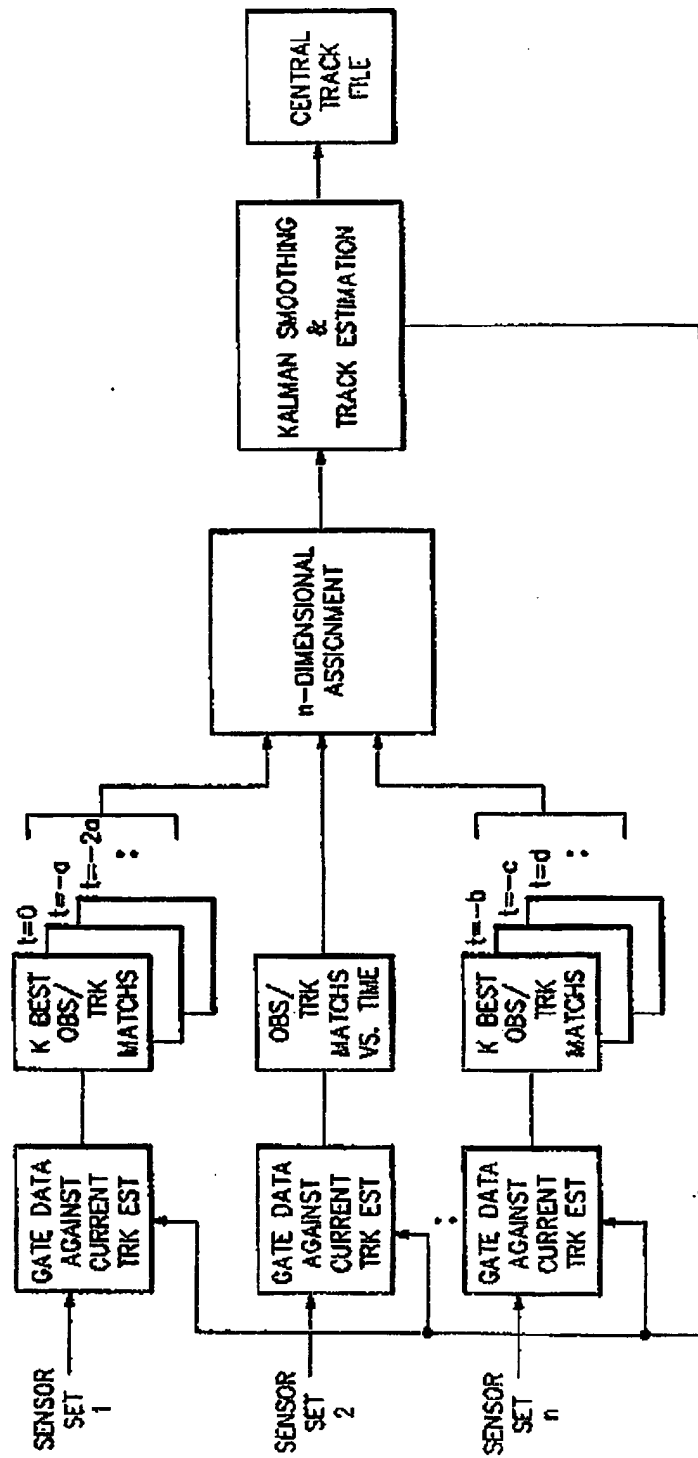


FIG. 24

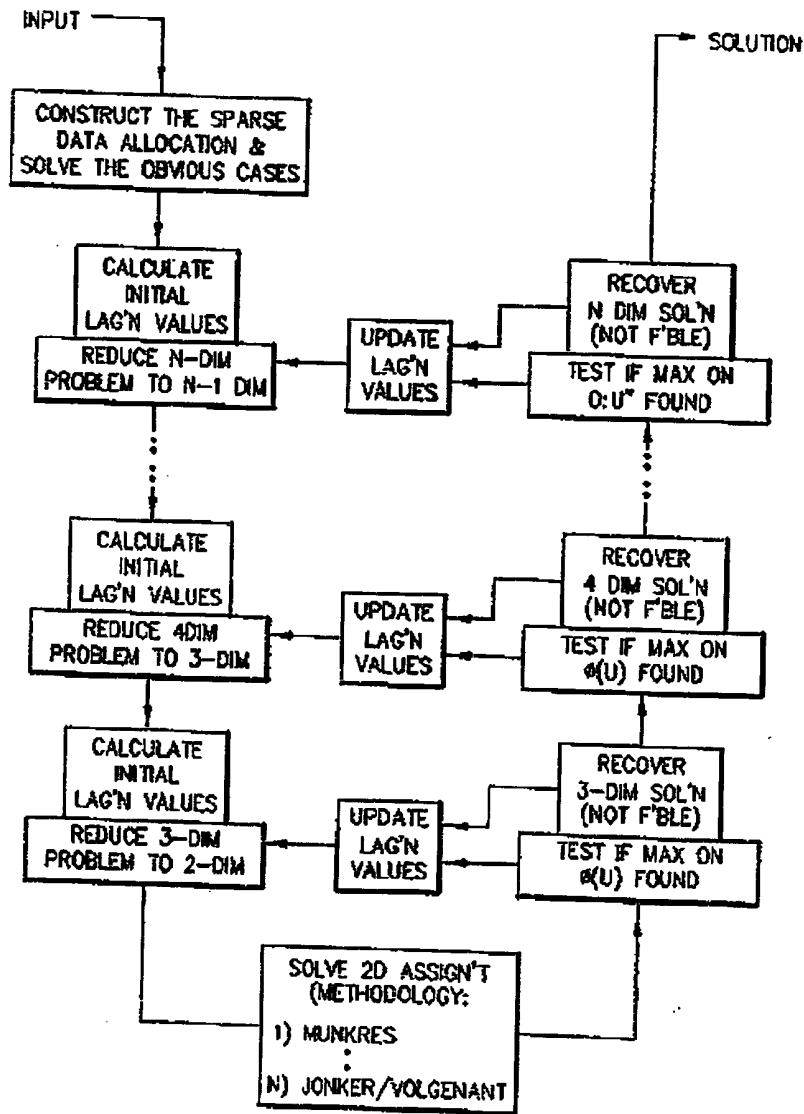
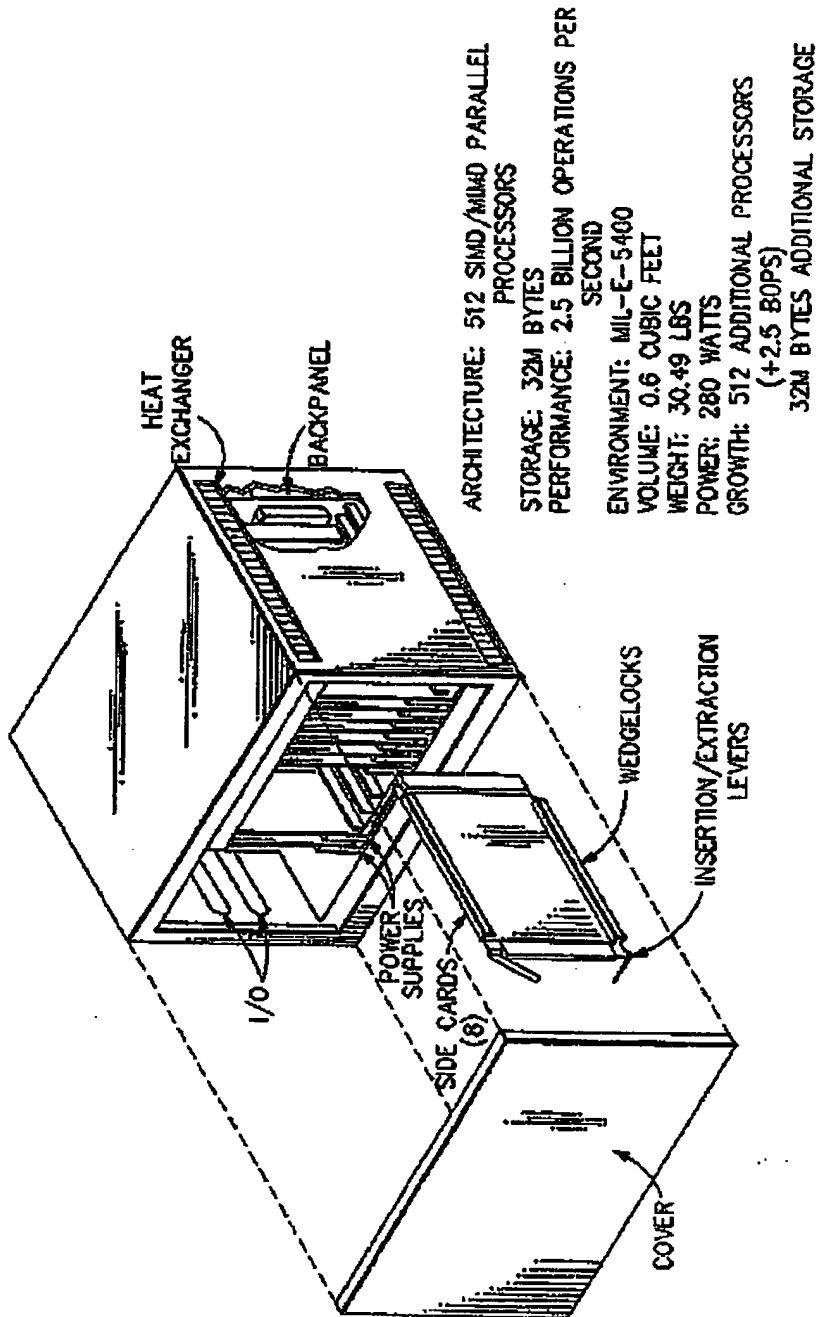
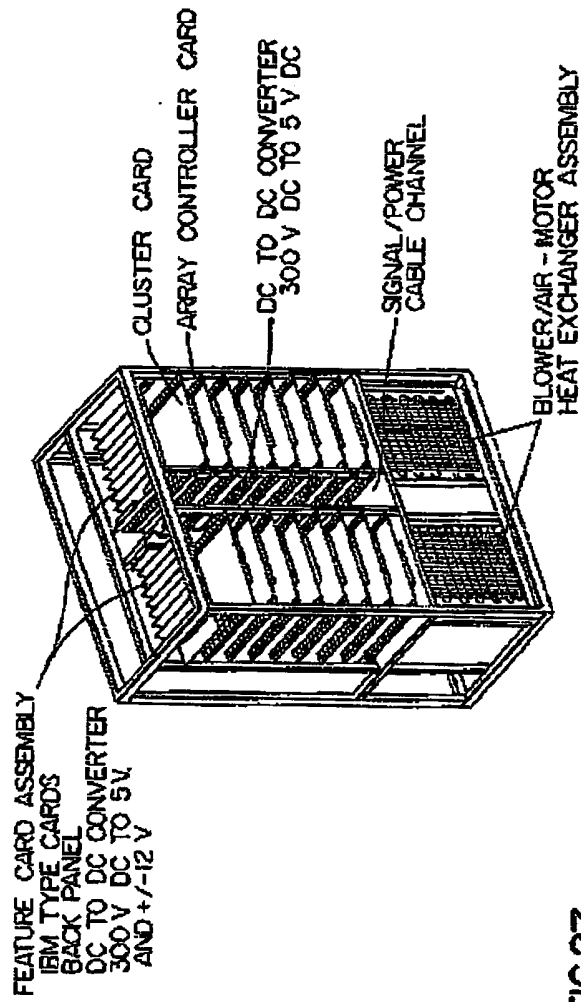


FIG.25



**FIG.26**



**FIG. 27**